

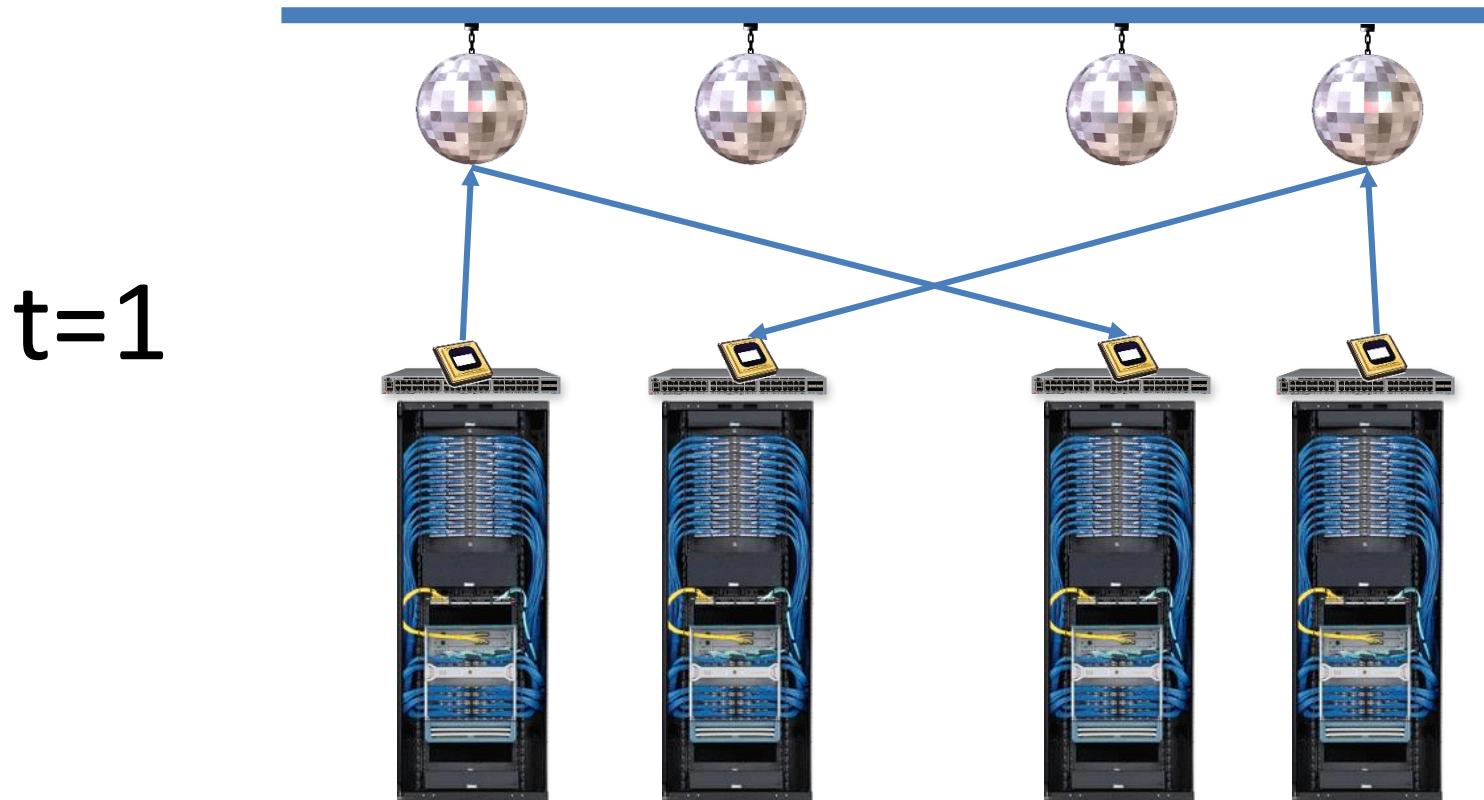
Networks in the Disco: Algorithms for Demand-Aware and Self-Adjusting Networks

Stefan Schmid (University of Vienna)



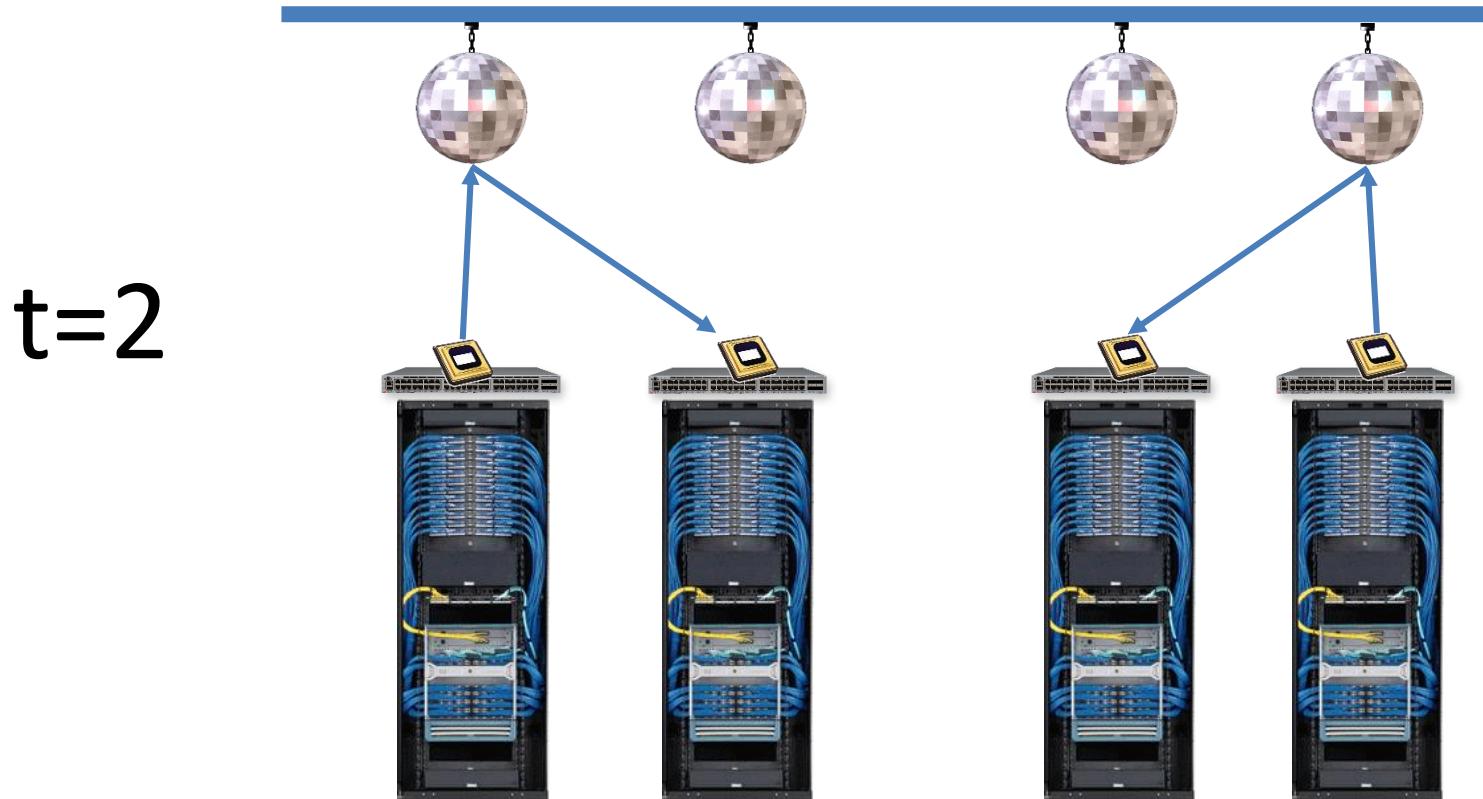
Motivation: Free-Space Optics

(*ProjectToR*)



Motivation: Free-Space Optics

(*ProjectToR*)

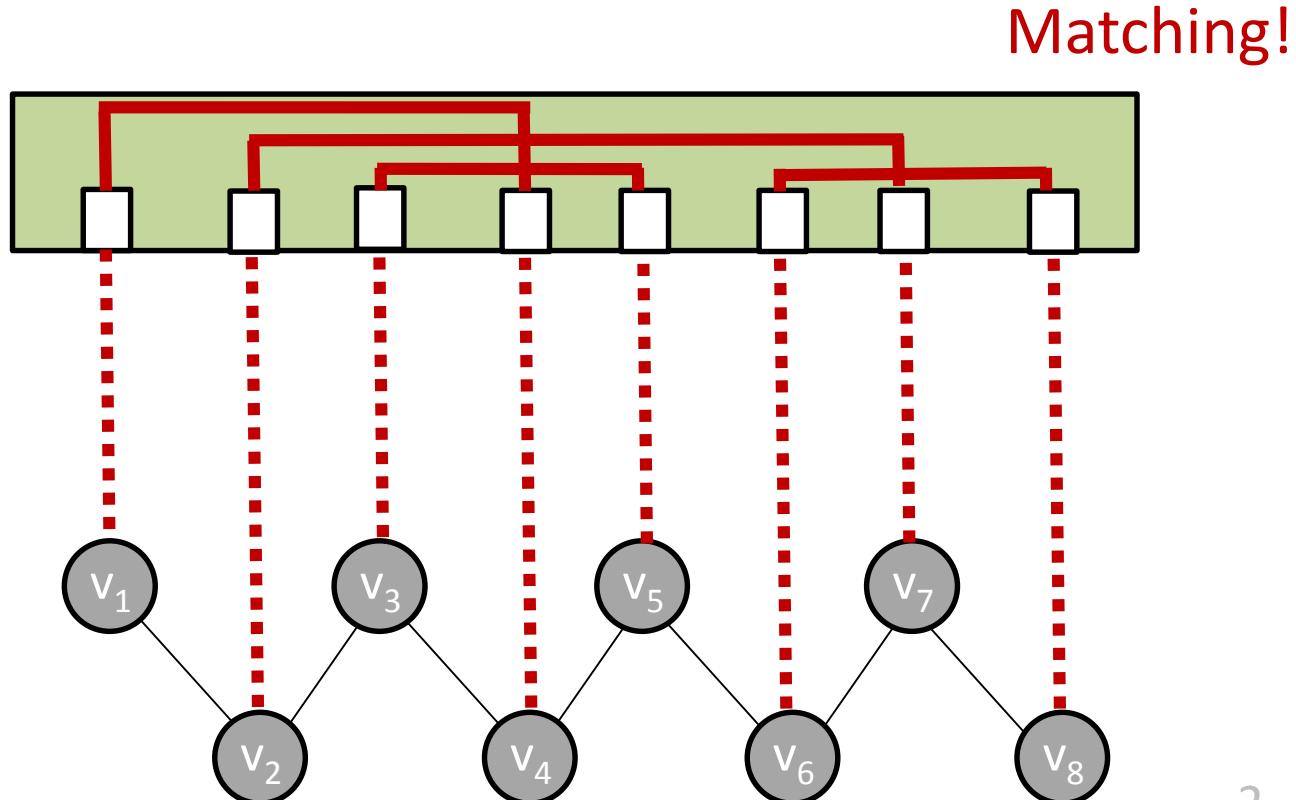


Also: Reconfigurable Optical Switches (*Helios*, *c-Through*, etc.)

Dynamic topology:
optical switch
(e.g. matching)

$t=1$

Static topology:
electric

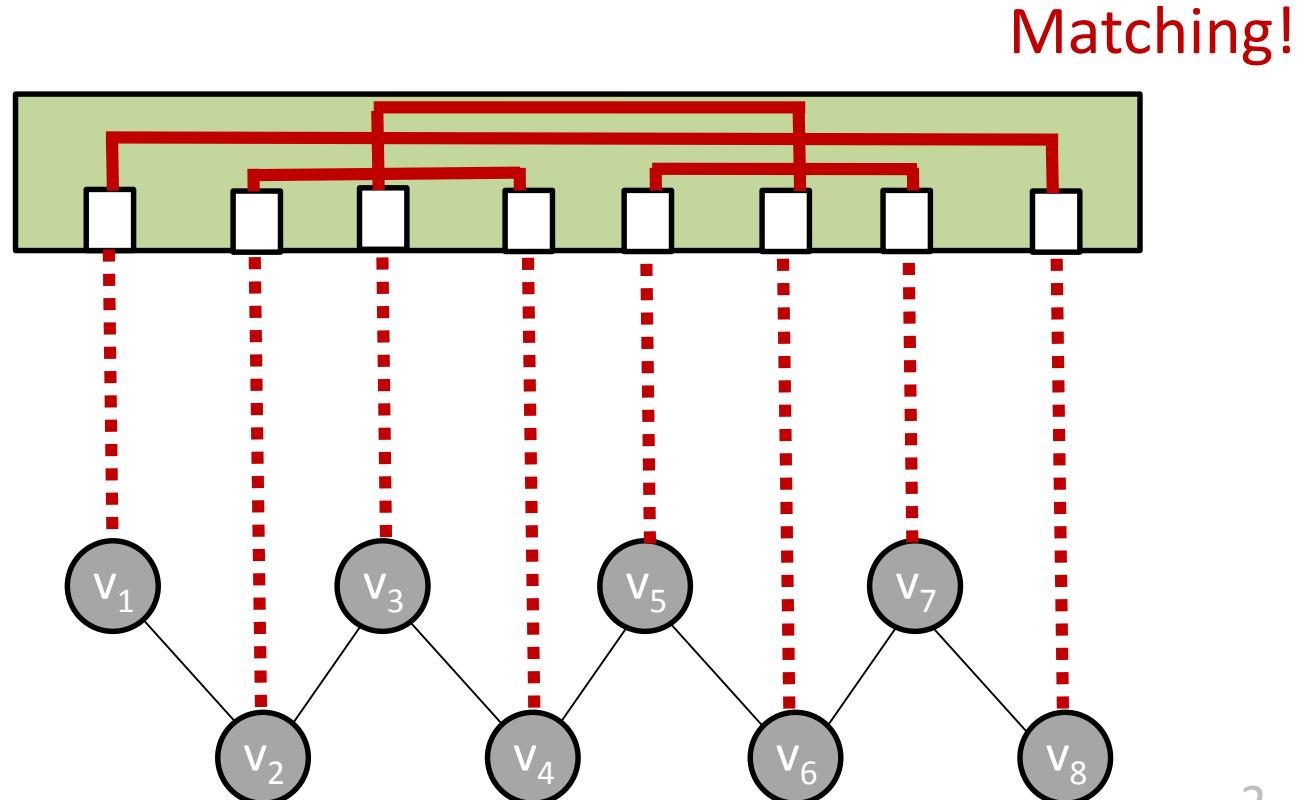


Also: Reconfigurable Optical Switches (*Helios*, *c-Through*, etc.)

Dynamic topology:
optical switch
(e.g. matching)

$t=2$

Static topology:
electric



Emerging Technologies

Movable Antennas

- Halperin et al. "Augmenting data center networks with multi-gigabit wireless links," SIGCOMM 2011.

60GHz Wireless Communication

- Zhou et al. "Mirror mirror on the ceiling: Flexible wireless links for data centers," CCR 2012.
- Kandula et al. "*Flyways* to de-congest data center networks," 2009.

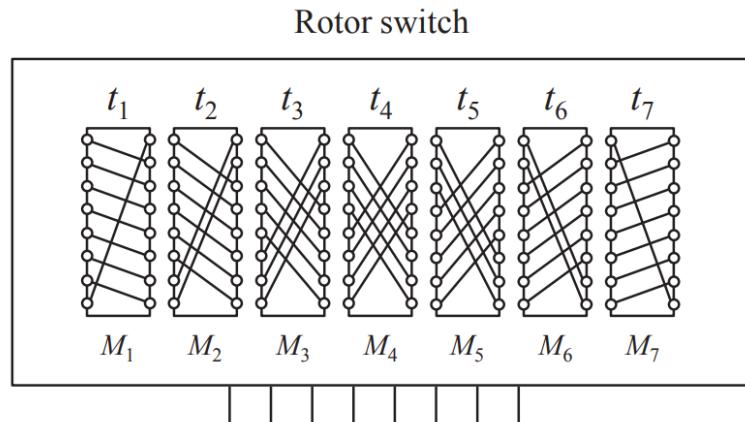
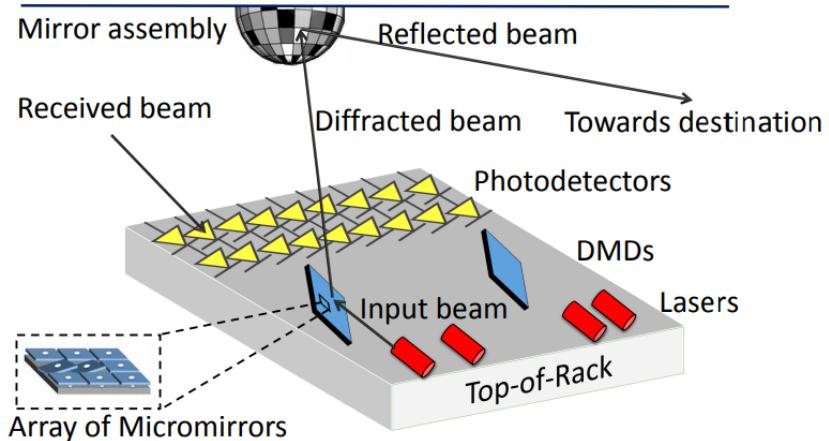
Free-Space Optics

- Ghobadi et al., "*Projector*: Agile reconfigurable data center interconnect," SIGCOMM 2016.
- Hamedazimi et al. "*Firefly*: A reconfigurable wireless data center fabric using free-space optics," CCR 2014.

Optical Circuit Switches

- Farrington et al. "*Helios*: a hybrid electrical/optical switch architecture for modular data centers," CCR 2010.
- Mellette et al. "*Rotornet*: A scalable, low-complexity, optical datacenter network," SIGCOMM 2017.
- Farrington et al. "Integrating microsecond circuit switching into the data center," SIGCOMM 2013.
- Liu et al. "Circuit switching under the radar with reactor.," NSDI 2014

Etc.!

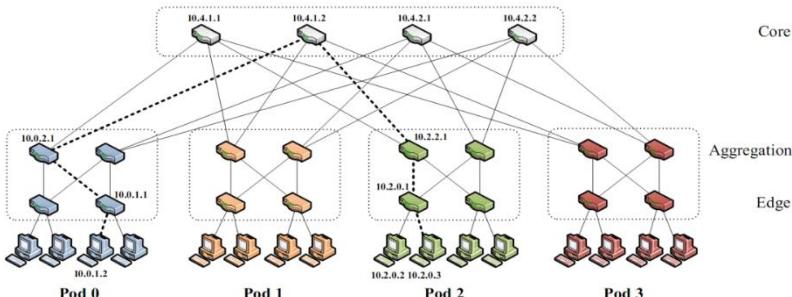




Observation: Technology Enables
“Demand-Aware Networks”

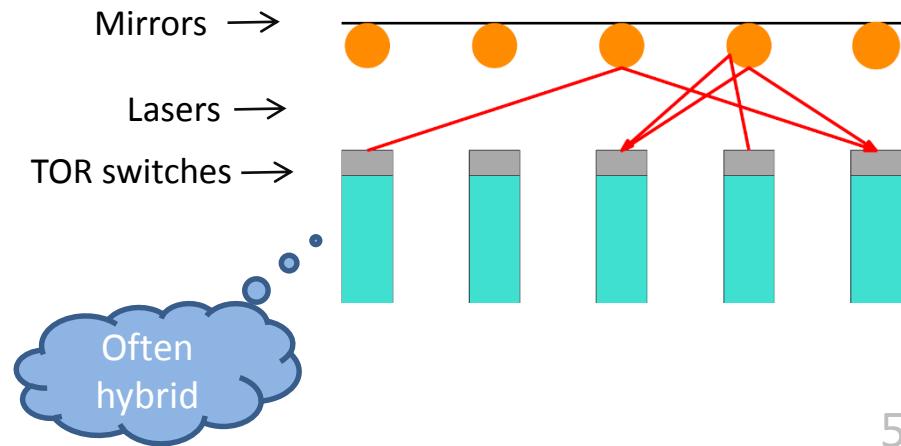
Traditional Networks

- Usually optimized **for the “worst-case”** (all-to-all communication)
- Example, fat-tree topologies: provide **full bisection bandwidth**
- Lower bounds and hard **trade-offs**, e.g., *degree vs diameter*



DANs

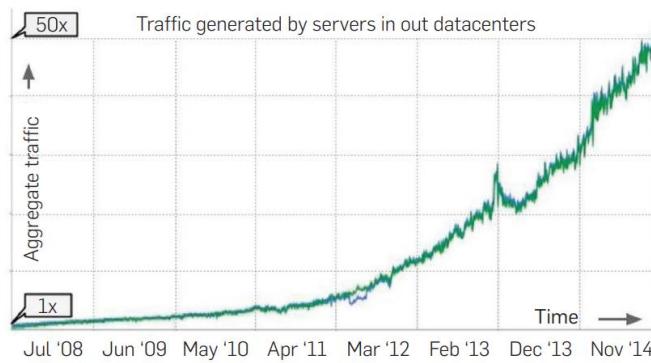
- **Demand-Aware Network (DAN)**
 - Optimized **toward the workload** it serves (e.g., **route length**)
 - Statically or **even dynamically**



Why...?

Growing Traffic and Cost...

Batch processing, web services,
distributed ML, ...: *data-centric applications* are distributed and interconnecting network is *critical*



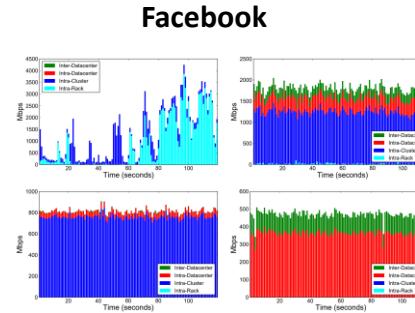
Source: Jupiter Rising. SIGCOMM 2015.

Aggregate server traffic in
Google's datacenter fleet

... But Much Structure!

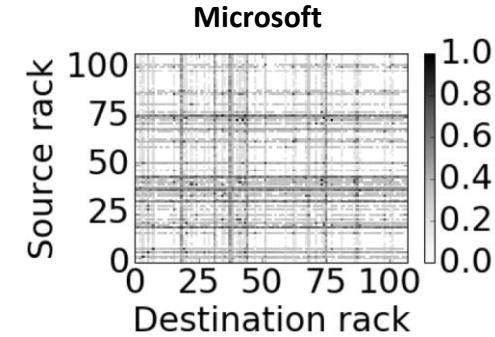
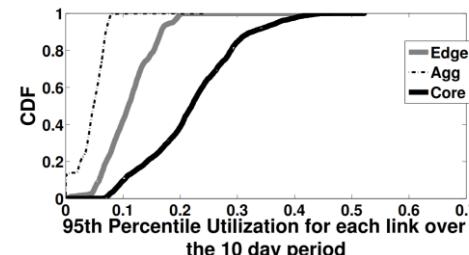


Spatial (*sparse!*) and temporal *locality*



Inside the Social Network's
(Datacenter) Network @
SIGCOMM 2015

Benson et al.

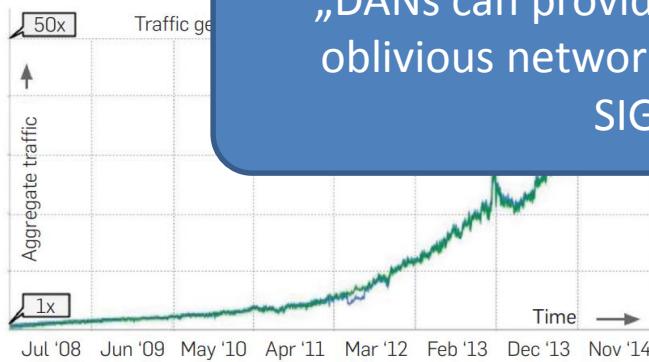


ProjecToR @ SIGCOMM 2016

Understanding Data Center Traffic
Characteristics @ WREN 2009

Growing Traffic and Cost...

Batch processing, web services, **distributed ML**, ...: ***data-centric applications*** are distributed and interconnecting network is ***critical***

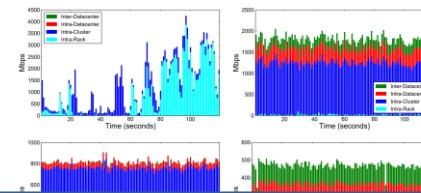


Source: Jupiter Rising. SIGCOMM 2015.

Aggregate server traffic in
Google's datacenter fleet

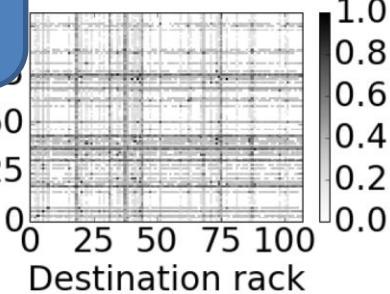
... But Much Structure!

Facebook



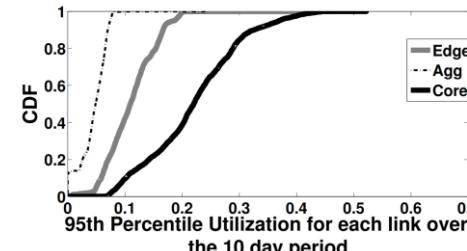
Spatial (***sparse!***) and temporal ***locality***

Microsoft



ProjecToR @ SIGCOMM 2016

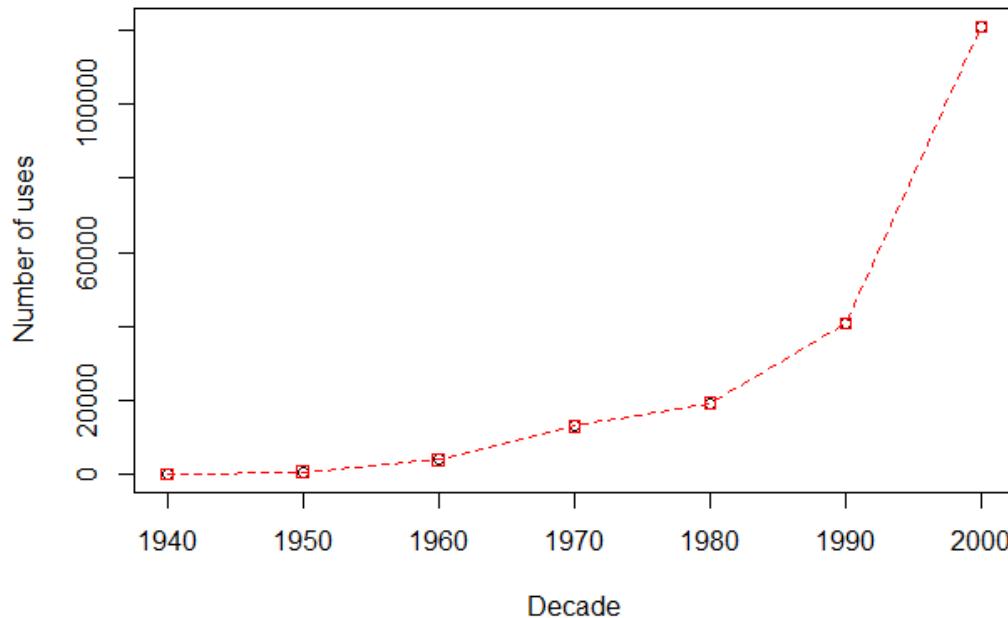
Benson et al.



Understanding Data Center Traffic Characteristics @ WREN 2009

Fun Fact

Use of the phrase 'Exponential growth' by decade



Data from Google Scholar

Roadmap

- Motivation: Demand-Aware Networks
- **Principles of Static Demand-Aware Network Designs**
- Principles of Dynamic Demand-Aware Network Designs
- Principles of Decentralized Approaches



A “Simple” DAN Design Problem

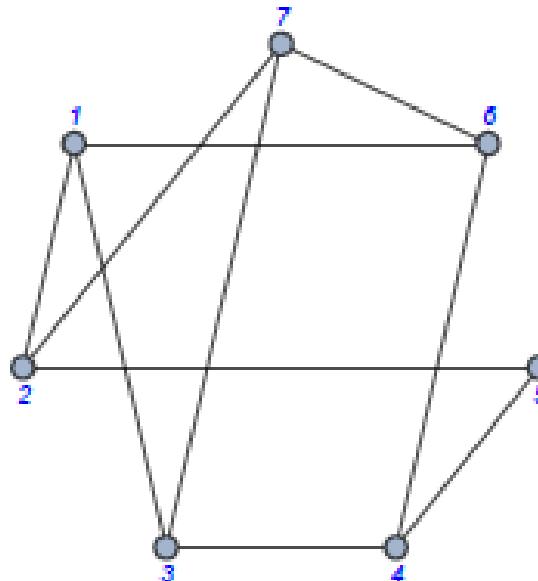
Input: Workload

Destinations

Sources	1	2	3	4	5	6	7
1	0	$\frac{2}{65}$	$\frac{1}{13}$	$\frac{1}{65}$	$\frac{1}{65}$	$\frac{2}{65}$	$\frac{3}{65}$
2	$\frac{2}{65}$	0	$\frac{1}{65}$	0	0	0	$\frac{2}{65}$
3	$\frac{1}{13}$	$\frac{1}{65}$	0	$\frac{2}{65}$	0	0	$\frac{1}{13}$
4	$\frac{1}{65}$	0	$\frac{2}{65}$	0	$\frac{4}{65}$	0	0
5	$\frac{1}{65}$	0	$\frac{3}{65}$	$\frac{4}{65}$	0	0	0
6	$\frac{2}{65}$	0	0	0	0	0	$\frac{3}{65}$
7	$\frac{3}{65}$	$\frac{2}{65}$	$\frac{1}{13}$	0	0	$\frac{3}{65}$	0

design

Output: DAN



Demand matrix: joint distribution

... of *constant degree* (scalability)

A “Simple” DAN Design Problem

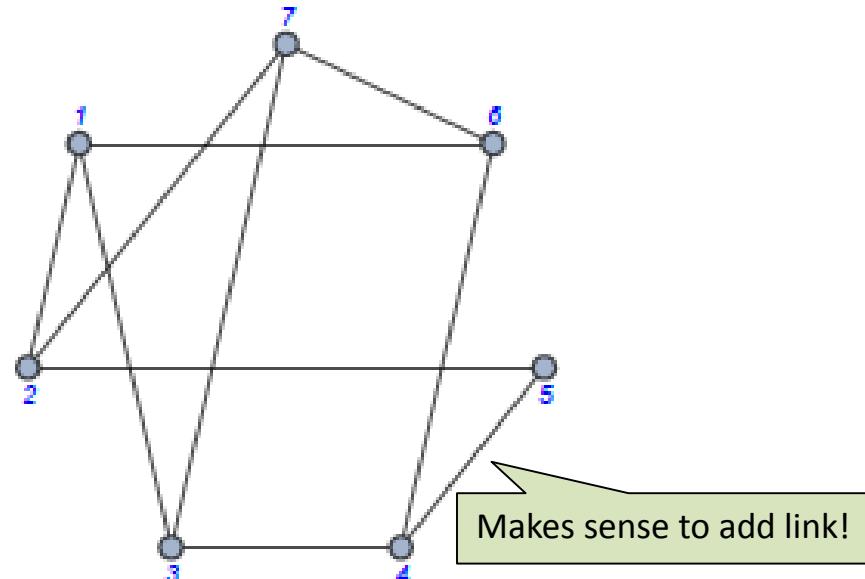
Input: Workload

Destinations

Sources	1	2	3	4	5	6	7
1	0	$\frac{2}{65}$	$\frac{1}{13}$	$\frac{1}{65}$	$\frac{1}{65}$	$\frac{2}{65}$	$\frac{3}{65}$
2	$\frac{2}{65}$	0	$\frac{1}{65}$	0	0	0	$\frac{2}{65}$
3	$\frac{1}{13}$	$\frac{1}{65}$	0	$\frac{2}{65}$	0	$\frac{13}{65}$	0
4	$\frac{1}{65}$	0	$\frac{2}{65}$	0	$\frac{4}{65}$	0	0
5	$\frac{1}{65}$	0	$\frac{3}{65}$	$\frac{4}{65}$	0	0	0
6	$\frac{2}{65}$	0	0	0	0	0	$\frac{3}{65}$
7	$\frac{3}{65}$	$\frac{2}{65}$	$\frac{1}{13}$	0	0	$\frac{3}{65}$	0

design

Output: DAN



Demand matrix: joint distribution

... of *constant degree* (scalability)

A “Simple” DAN Design Problem

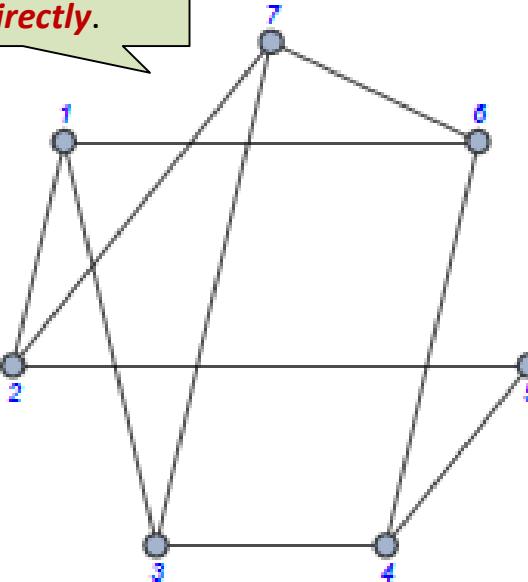
Input: Workload

Destinations								
		1 communicates to many.						
		6	7					
1	0	$\frac{2}{65}$	$\frac{1}{13}$	$\frac{1}{65}$	$\frac{1}{65}$	$\frac{2}{65}$	$\frac{3}{65}$	
2	$\frac{2}{65}$	0	$\frac{1}{65}$	0	0	0	$\frac{2}{65}$	
3	$\frac{1}{13}$	$\frac{1}{65}$	0	$\frac{2}{65}$	0	0	$\frac{1}{13}$	
4	$\frac{1}{65}$	0	$\frac{2}{65}$	0	$\frac{4}{65}$	0	0	
5	$\frac{1}{65}$	0	$\frac{3}{65}$	$\frac{4}{65}$	0	0	0	
6	$\frac{2}{65}$	0	0	0	0	0	$\frac{3}{65}$	
7	$\frac{3}{65}$	$\frac{2}{65}$	$\frac{1}{13}$	0	0	$\frac{3}{65}$	0	

Output: DAN

Bounded degree: route
to 7 *indirectly*.

design



Demand matrix: joint distribution

... of *constant degree* (scalability)

A “Simple” DAN Design Problem

Input: Workload

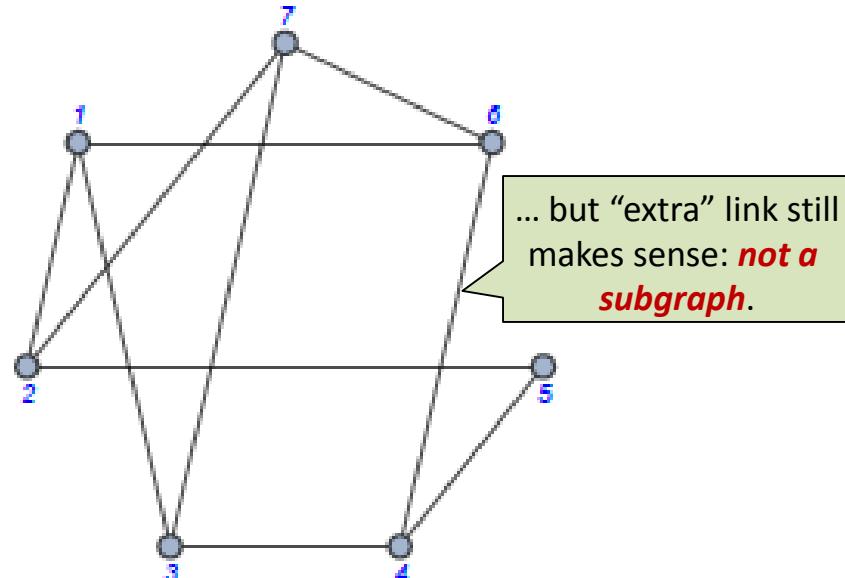
Destinations

Sources	1	2	3	4	5	6	7
1	0	$\frac{2}{65}$	$\frac{1}{13}$	$\frac{1}{65}$	$\frac{1}{65}$	$\frac{2}{65}$	$\frac{3}{65}$
2	$\frac{2}{65}$	0	$\frac{1}{65}$	0	0	0	$\frac{2}{65}$
3	$\frac{1}{13}$	$\frac{1}{65}$	0	$\frac{2}{65}$	0	0	$\frac{1}{13}$
4	$\frac{1}{65}$	0	$\frac{2}{65}$	0	$\frac{4}{65}$	0	0
5	$\frac{1}{65}$	0	$\frac{3}{65}$	$\frac{4}{65}$	0		
6	$\frac{2}{65}$	0	0	0	0		
7	$\frac{3}{65}$	$\frac{2}{65}$	$\frac{1}{13}$	0	0	$\frac{3}{65}$	0

design

4 and 6 don't communicate...

Output: DAN



Demand matrix: joint distribution

... of *constant degree* (scalability)

Case Study: DAN for Short Routes

Shorter routes: smaller bandwidth footprint, lower latency, less energy, ...

More Formally: DAN Design Problem

Input:

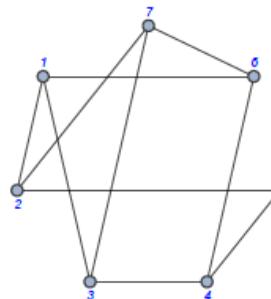
$\mathcal{D}[p(i,j)]$: joint **distribution**, Δ

		Y						
		1	2	3	4	5	6	7
X	1	0	<u>2</u>	<u>1</u>	<u>1</u>	<u>1</u>	<u>2</u>	<u>3</u>
	2	<u>2</u>	0	<u>1</u>	0	0	0	<u>2</u>
3	<u>1</u>	<u>1</u>	0	<u>2</u>	0	0	0	<u>1</u>
	<u>13</u>	<u>65</u>	<u>65</u>	<u>65</u>	<u>65</u>	<u>65</u>	<u>65</u>	<u>13</u>
4	<u>1</u>	0	<u>2</u>	0	<u>4</u>	0	0	<u>0</u>
	<u>65</u>	<u>0</u>						
5	<u>1</u>	0	<u>3</u>	<u>4</u>	0	0	0	<u>0</u>
	<u>65</u>	<u>0</u>						
6	<u>2</u>	0	0	0	0	0	0	<u>65</u>
	<u>65</u>	<u>3</u>						
7	<u>3</u>	<u>2</u>	<u>1</u>	0	0	<u>3</u>	0	<u>0</u>
	<u>65</u>	<u>65</u>	<u>13</u>	<u>0</u>	<u>0</u>	<u>65</u>	<u>0</u>	<u>0</u>



Output:

N: DAN



Bounded degree
 $\Delta=3$

Objective:

Expected Path Length (EPL):

Demand-weighted route length

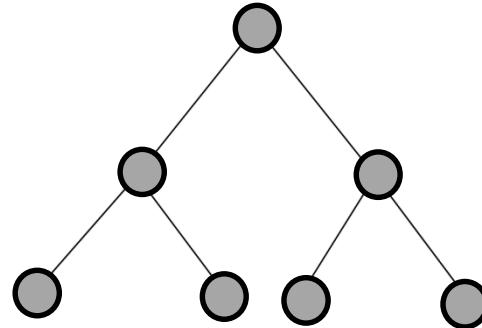
$$\text{EPL}(\mathcal{D}, N) = \sum_{(u,v) \in \mathcal{D}} p(u, v) \cdot d_N(u, v)$$

Path length **on DAN N.**

Frequency

Some Examples

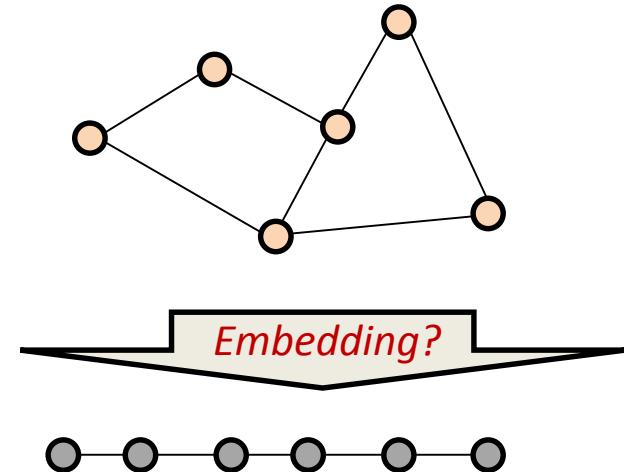
- DANs of $\Delta = 3$:
 - E.g., complete binary **tree**
 - $d_N(u,v) \leq 2 \log n$
 - Can we do **better** than **$\log n$** ?
- DANs of $\Delta = 2$:
 - E.g., set of **lines** and **cycles**



How hard is it to design a DAN?

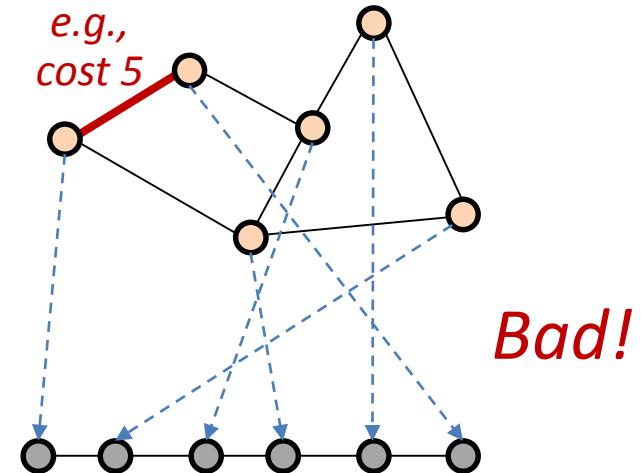
DAN design can be NP-hard

- **Example $\Delta = 2$:** A Minimum Linear Arrangement (**MLA**) problem
 - A “Virtual Network Embedding Problem”, VNEP
 - *Minimize sum* of lengths of virtual edges



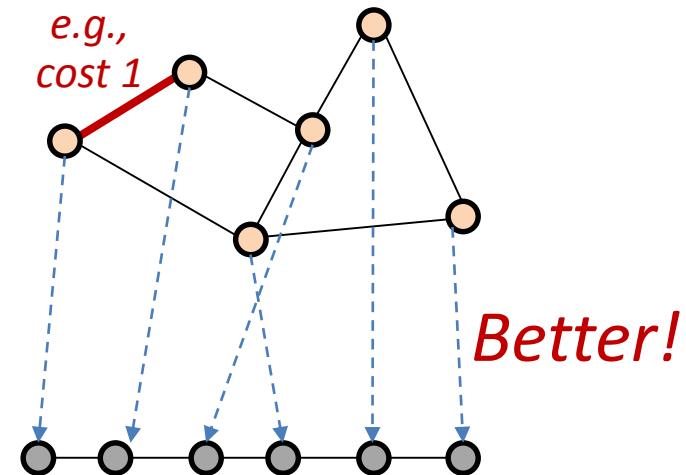
DAN design can be NP-hard

- **Example $\Delta = 2$: A Minimum Linear Arrangement (MLA) problem**
 - A “Virtual Network Embedding Problem”, VNEP
 - *Minimize sum* of lengths of virtual edges



DAN design can be NP-hard

- **Example $\Delta = 2$:** A Minimum Linear Arrangement (**MLA**) problem
 - A “Virtual Network Embedding Problem”, VNEP
 - *Minimize sum* of lengths of virtual edges

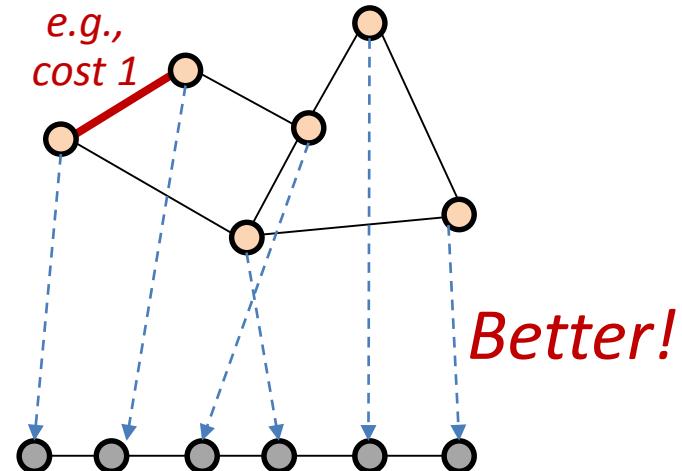


Better!

DAN design can be NP-hard

- Example $\Delta = 2$: A Minimum Linear Arrangement (MLA) problem
 - A “Virtual Network Embedding Problem”, VNEP
 - *Minimizing the lengths of virtual edges*

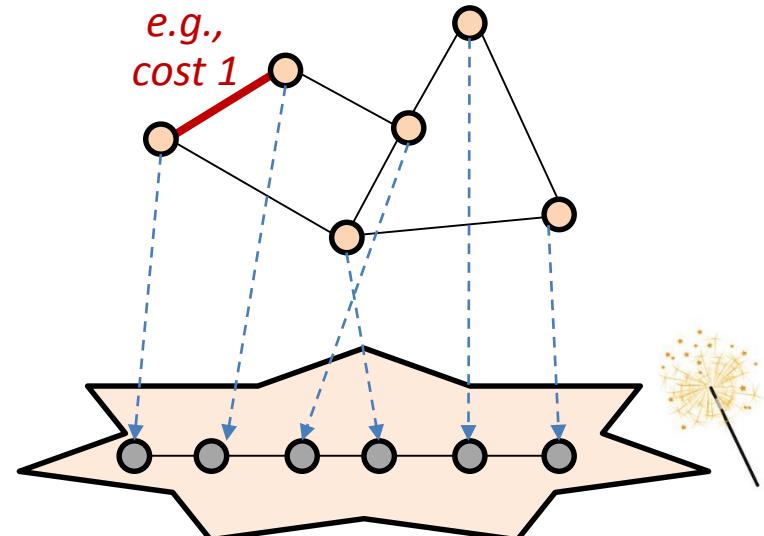
NP-hard, and so is DAN design.



DAN design can be NP-hard

- Example $\Delta = 2$: A Minimum Linear Arrangement (MLA) problem
 - A “Virtual Network Embedding Problem”, VNEP
 - *Minimizing lengths of virtual edges*
- But what about > 2 ? *Embedding* problem still hard, but we have an additional **degree of freedom**:

Do topological flexibilities make problem easier or harder?!



A new knob for optimization!

So: How useful are DANs?

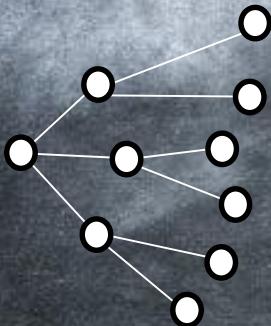
As always in computer science (e.g., also in coding, in self-adjusting datastructures, etc.): ***it depends!*** ☺

Expected Path Length in Traditional Networks?

Theorem (Traditional Networks):

Constant-degree networks have at least logarithmic diameter.

Proof.



1 Δ $\Delta(\Delta-1)$

In k steps, reach at most $1 + \sum \Delta(\Delta - 1)^i$

Each network with n nodes and max degree $\Delta > 2$ must have a diameter of at least $\log(n)/\log(\Delta-1)-1$.



Example: Clos, Bcube, Xpander.

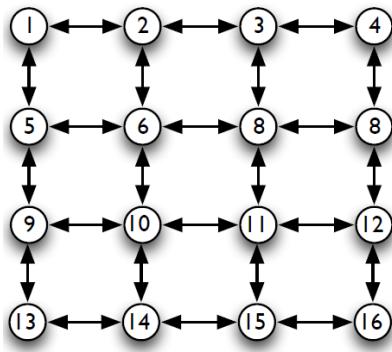
Can DANs do better?

Can DANs do better?

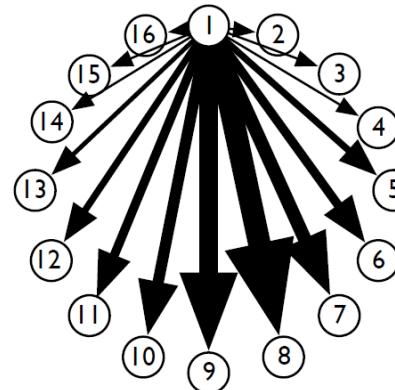
In general not really, e.g. in all-to-all communication
(*clique*): ***logarithmic diameter unavoidable.***

But sometimes, DANs can be much better!

Example 1: low-degree demand



Example 2: high-degree but *skewed* demand



- Already low degree: degree-4 DAN can serve this *at cost 1*.
- If sufficiently skewed: constant-degree DAN can serve it at cost ***O(1)***

So on what does it depend?

So on what does it depend?



We argue: on the
“entropy” of the demand!

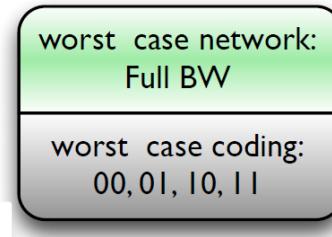




00110101...



if demand **arbitrary** and **unknown**



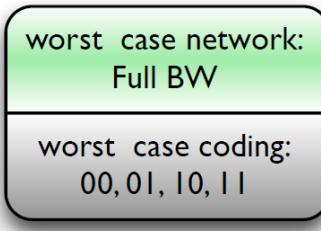
log diameter

log # bits / symbol



An Analogy to Coding

if demand **arbitrary** and **unknown**



log diameter

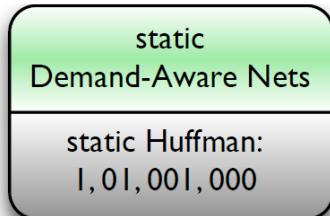
log # bits / symbol



DAN!

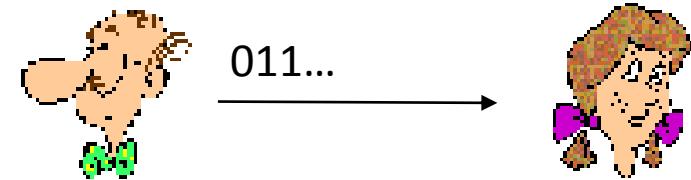


if demand **known** and **fixed**

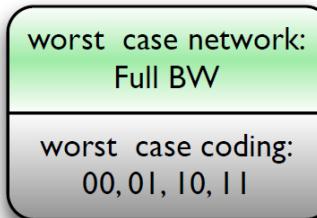


entropy?

entropy / symbol



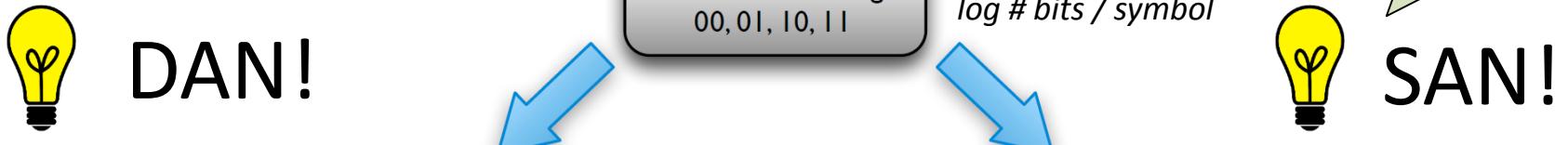
if demand **arbitrary** and **unknown**



log diameter

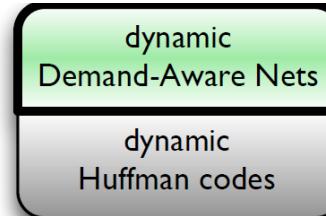
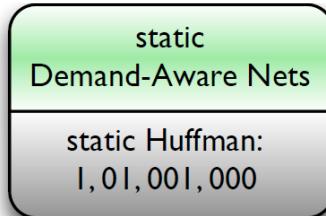
Dynamic DANs:
Aka. **Self-Adjusting Networks (SANs)!**

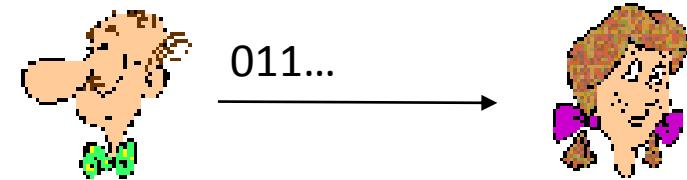
log # bits / symbol



if demand **known** and **fixed** if demand **unknown** but **reconfigurable**

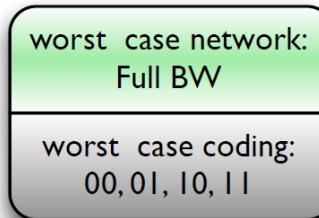
entropy?
entropy / symbol





An Analogy to Coding

if demand **arbitrary** and **unknown**



log diameter

log # bits / symbol

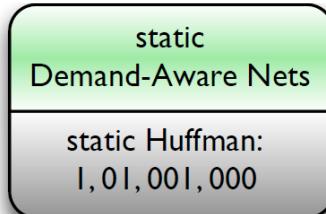
Dynamic DANs:
Aka. **Self-Adjusting Networks (SANs)!**



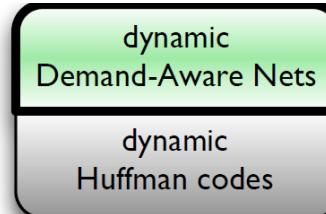
DAN!



if demand **known** and **fixed**



Can exploit
spatial locality!



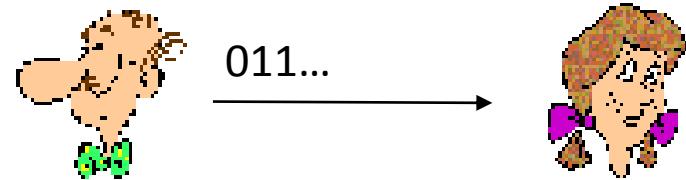
SAN!



if demand **unknown** but **reconfigurable**

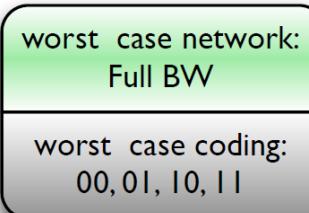


Additionally exploit
temporal locality!



An Analogy to Coding

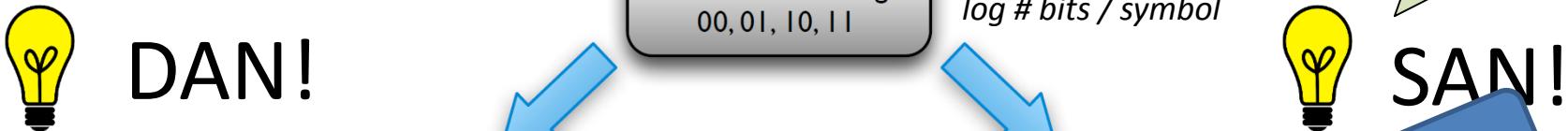
if demand **arbitrary** and **unknown**



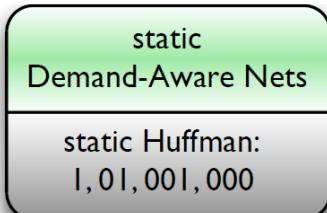
log diameter

Dynamic DANs:
Aka. **Self-Adjusting Networks (SANs)!**

log # bits / symbol

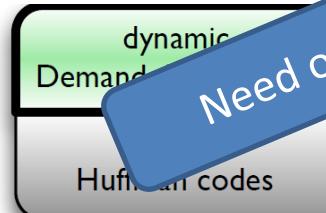


if demand **known** and **fixed**



Can exploit
spatial locality!

if demand **unknown** but **repeating**



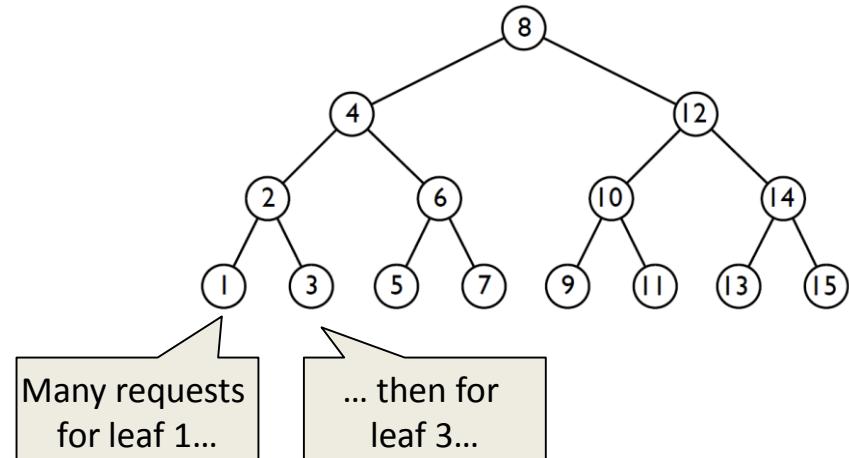
Need online algorithms!



Additionally exploit
temporal locality!

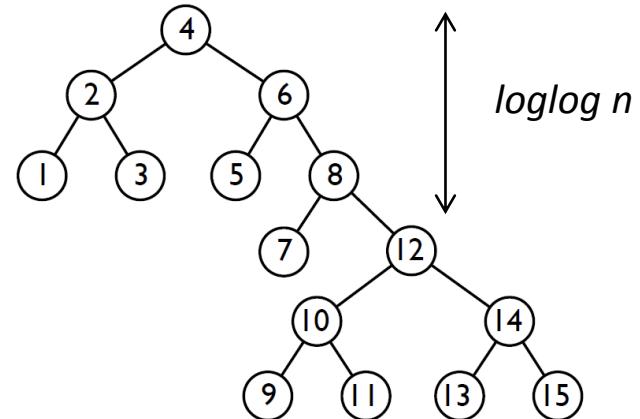
Analogous to *Datastructures*: Oblivious...

- Traditional, **fixed** BSTs do not rely on any assumptions on the demand
- Optimize for the **worst-case**
- Example **demand**:
 $1, \dots, 1, 3, \dots, 3, 5, \dots, 5, 7, \dots, 7, \dots, \log(n), \dots, \log(n)$
 $\longleftrightarrow \longleftrightarrow \longleftrightarrow \longleftrightarrow \longleftrightarrow \quad \longleftrightarrow$
many many many many *many*
- Items stored at **$O(\log n)$** from the root, **uniformly** and **independently** of their frequency



... Demand-Aware ...

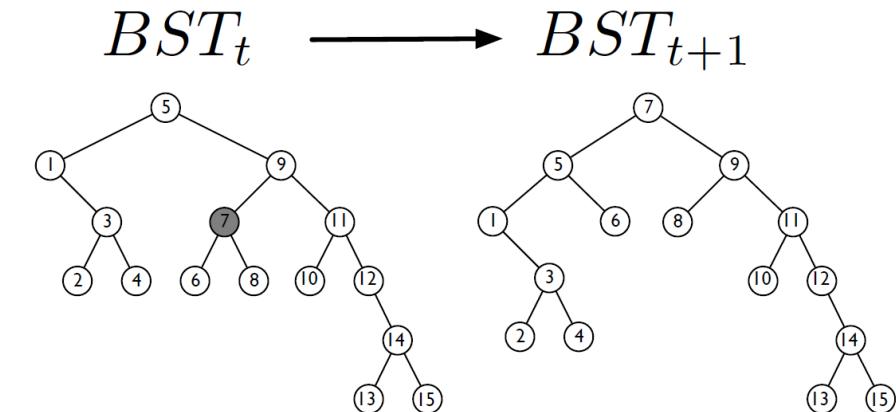
- **Demand-aware fixed** BSTs can take advantage of *spatial locality* of the demand
- E.g.: place frequently accessed elements close to the root
- E.g., **Knuth/Mehlhorn/Tarjan** trees
- Recall example **demand**:
 $1, \dots, 1, 3, \dots, 3, 5, \dots, 5, 7, \dots, 7, \dots, \log(n), \dots, \log(n)$
 - Amortized cost $O(\log \log n)$



 Amortized cost corresponds
to *empirical entropy of demand!*

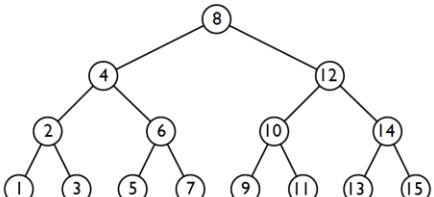
... Self-Adjusting!

- Demand-aware reconfigurable BSTs can additionally take advantage of *temporal locality*
- By moving accessed element to the root: amortized cost is *constant*, i.e., $O(1)$
 - Recall example **demand**:
 $1, \dots, 1, 3, \dots, 3, 5, \dots, 5, 7, \dots, 7, \dots, \log(n), \dots, \log(n)$

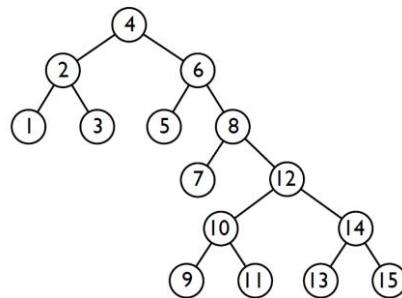


Datastructures

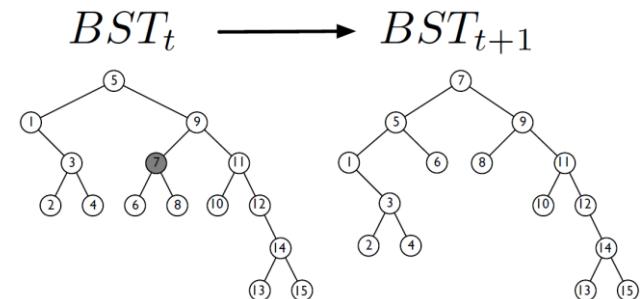
Oblivious



Demand-Aware



Self-Adjusting



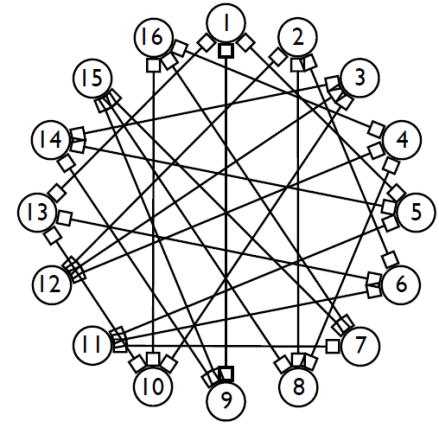
Lookup
 $O(\log n)$

Exploit **spatial locality**:
empirical entropy $O(\log \log n)$

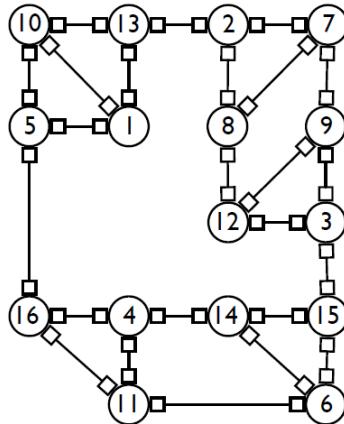
Exploit **temporal locality** as well:
 $O(1)$

Analogously for Networks

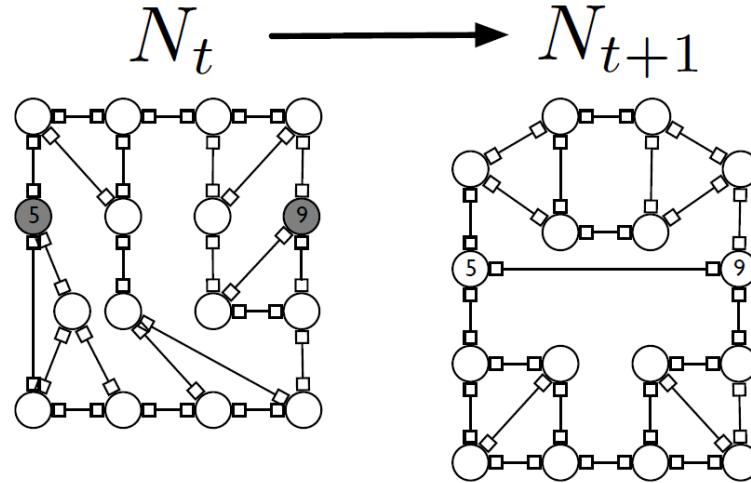
Oblivious



DAN



SAN



Const degree
(e.g., **expander**):
route lengths **$O(\log n)$**

Exploit **spatial locality**

Exploit **temporal locality** as well

SPEED
LIMIT

?

Limitations of (Static) DANs: Entropy-Based Lower Bounds



Lower Bound Idea: Leverage Coding or Datastructure!

Destinations		1	2	3	4	5	6	7
Sources	1	0 65	$\frac{2}{65}$	$\frac{1}{13}$	$\frac{1}{65}$	$\frac{1}{65}$	$\frac{2}{65}$	$\frac{3}{65}$
	2	$\frac{2}{65}$	0 65	$\frac{1}{65}$	0 0	0 0	0 0	$\frac{2}{65}$
	3	$\frac{1}{13}$	$\frac{1}{65}$	0 65	$\frac{2}{65}$	0 0	0 0	$\frac{1}{13}$
	4	$\frac{1}{65}$	0 65	$\frac{2}{65}$	0 65	$\frac{4}{65}$	0 0	0 0
	5	$\frac{1}{65}$	0 65	$\frac{3}{65}$	$\frac{4}{65}$	0 0	0 0	0 0
	6	$\frac{2}{65}$	0 65	0 0	0 0	0 0	0 0	$\frac{3}{65}$
	7	$\frac{3}{65}$	$\frac{2}{65}$	$\frac{1}{13}$	0 0	0 0	$\frac{3}{65}$	0 0

- DAN just for a **single (source) node 1**: cannot do better than Δ -ary **Huffman tree** for its destinations [0,1/65,1/13,1/65,1/65,2/65,3/65]
 - resp. **Knuth/Mehlhorn/Tarjan** tree if search property required
- How good can this tree be?



Entropy lower bound on EPL known for binary trees, e.g. **Mehlhorn** 1975 for BST

Lower Bound Idea: Leverage Coding or Data Structures

Destinations		1	2	3	4	5	6	7
Sources	1	0 65	$\frac{2}{65}$	$\frac{1}{13}$	$\frac{1}{65}$	$\frac{1}{65}$	$\frac{2}{65}$	$\frac{3}{65}$
	2	$\frac{2}{65}$	0 65	$\frac{1}{65}$	0 65	0 65	0 65	$\frac{2}{65}$
	3	$\frac{1}{13}$	$\frac{1}{65}$	0 65	$\frac{2}{65}$	0 65	0 65	$\frac{1}{13}$
	4	$\frac{1}{65}$	0 65	$\frac{2}{65}$	0 65	$\frac{4}{65}$	0 65	0 65
	5	$\frac{1}{65}$	0 65	$\frac{3}{65}$	$\frac{4}{65}$	0 65	0 65	0 65
	6	$\frac{2}{65}$	0 65	0 65	0 65	0 65	0 65	$\frac{3}{65}$
	7	$\frac{3}{65}$	$\frac{2}{65}$	$\frac{1}{13}$	0 65	0 65	$\frac{3}{65}$	0 65

- DAN just for a **single (source) node 1**: cannot do better than Δ -ary **Huffman tree** for its destinations [0,1/65,1/13,1/65,1/65,2/65,3/65]
 - resp. **Knuth/Mehlhorn/Tarjan** tree if search property required

- How good can this tree be?



Entropy lower bound on EPL known for binary trees, e.g. **Mehlhorn** 1975 for BST

Lower Bound Idea: Leverage Coding or Data Structures

Destinations		1	2	3	4	5	6	7
Sources	1	0 65	$\frac{2}{65}$	$\frac{1}{13}$	$\frac{1}{65}$	$\frac{1}{65}$	$\frac{2}{65}$	$\frac{3}{65}$
	2	$\frac{2}{65}$	0 65	$\frac{1}{65}$	0 65	0 65	0 65	$\frac{2}{65}$
	3	$\frac{1}{13}$	$\frac{1}{65}$	0 65	$\frac{2}{65}$	0 65	0 65	$\frac{1}{13}$
	4	$\frac{1}{65}$	0 65	$\frac{2}{65}$	0 65	$\frac{4}{65}$	0 65	0 65
	5	$\frac{1}{65}$	0 65	$\frac{3}{65}$	$\frac{4}{65}$	0 65	0 65	0 65
	6	$\frac{2}{65}$	0 65	0 65	0 65	0 65	0 65	$\frac{3}{65}$
	7	$\frac{3}{65}$	$\frac{2}{65}$	$\frac{1}{13}$	0 65	0 65	$\frac{3}{65}$	0 65

An optimal “ego-tree”
for this source!

- DAN just for a **single (source) node 1**: cannot do better than Δ -ary **Huffman tree** for its destinations [0,1/65,1/13,1/65,1/65,2/65,3/65]
 - resp. **Knuth/Mehlhorn/Tarjan** tree if search property required
- How good can this tree be?



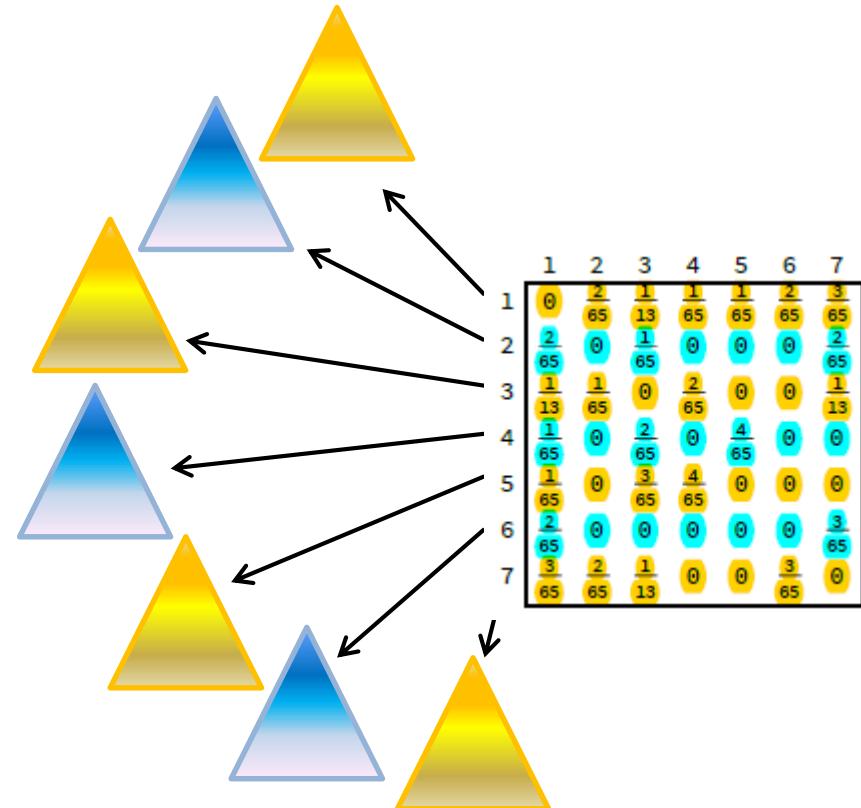
Entropy lower bound on EPL known for binary trees, e.g. **Mehlhorn** 1975 for BST



So: what is the entropy of the whole demand?

Lower Bound & Entropy of the Demand

- Proof idea ($EPL = \Omega(H_A(Y|X))$):
 - sources
 - destinations
 - entropy
- Compute **ego-tree** for each source node
- Take **union** of all **ego-trees**
- Violates **degree restriction** but valid lower bound



Lower Bound & Entropy of the Demand: Sources + Destinations

Do this in **both dimensions**:

$$\text{EPL} \geq \Omega(\max\{\mathcal{H}_\Delta(Y|X), \mathcal{H}_\Delta(X|Y)\})$$

$\Omega(\mathcal{H}_\Delta(X Y))$						
1	2	3	4	5	6	7
1	0 $\frac{2}{65}$	$\frac{1}{13}$	$\frac{1}{65}$	$\frac{1}{65}$	$\frac{2}{65}$	$\frac{3}{65}$
2	$\frac{2}{65}$	0 $\frac{1}{65}$	0	0	0 $\frac{2}{65}$	
3	$\frac{1}{13}$	$\frac{1}{65}$	0 $\frac{2}{65}$	0	0 $\frac{1}{13}$	
4	$\frac{1}{65}$	0 $\frac{2}{65}$	0	$\frac{4}{65}$	0 $\frac{3}{65}$	0
5	$\frac{1}{65}$	0 $\frac{3}{65}$	$\frac{4}{65}$	0	0 $\frac{3}{65}$	0
6	$\frac{2}{65}$	0 $\frac{3}{65}$	0	0	0 $\frac{3}{65}$	
7	$\frac{3}{65}$	$\frac{2}{65}$	$\frac{1}{13}$	0	0 $\frac{3}{65}$	0

\mathcal{D}

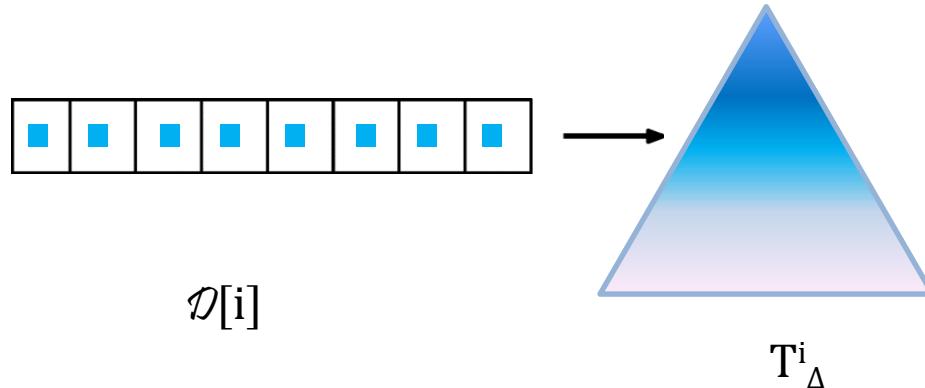
$\Omega(\mathcal{H}_\Delta(Y X))$						
1	2	3	4	5	6	7
1	0 $\frac{2}{65}$	$\frac{1}{13}$	$\frac{1}{65}$	$\frac{1}{65}$	$\frac{2}{65}$	$\frac{3}{65}$
2	$\frac{2}{65}$	0 $\frac{1}{65}$	0	0	0 $\frac{2}{65}$	$\frac{2}{65}$
3	$\frac{1}{13}$	$\frac{1}{65}$	0 $\frac{2}{65}$	0	0 $\frac{1}{13}$	
4	$\frac{1}{65}$	0 $\frac{2}{65}$	0	$\frac{4}{65}$	0 $\frac{3}{65}$	0
5	$\frac{1}{65}$	0 $\frac{3}{65}$	$\frac{4}{65}$	0	0 $\frac{3}{65}$	0
6	$\frac{2}{65}$	0 $\frac{3}{65}$	0	0	0 $\frac{3}{65}$	
7	$\frac{3}{65}$	$\frac{2}{65}$	$\frac{1}{13}$	0	0 $\frac{3}{65}$	0

Can DANs Match The Entropy Speed Limit? Upper Bounds



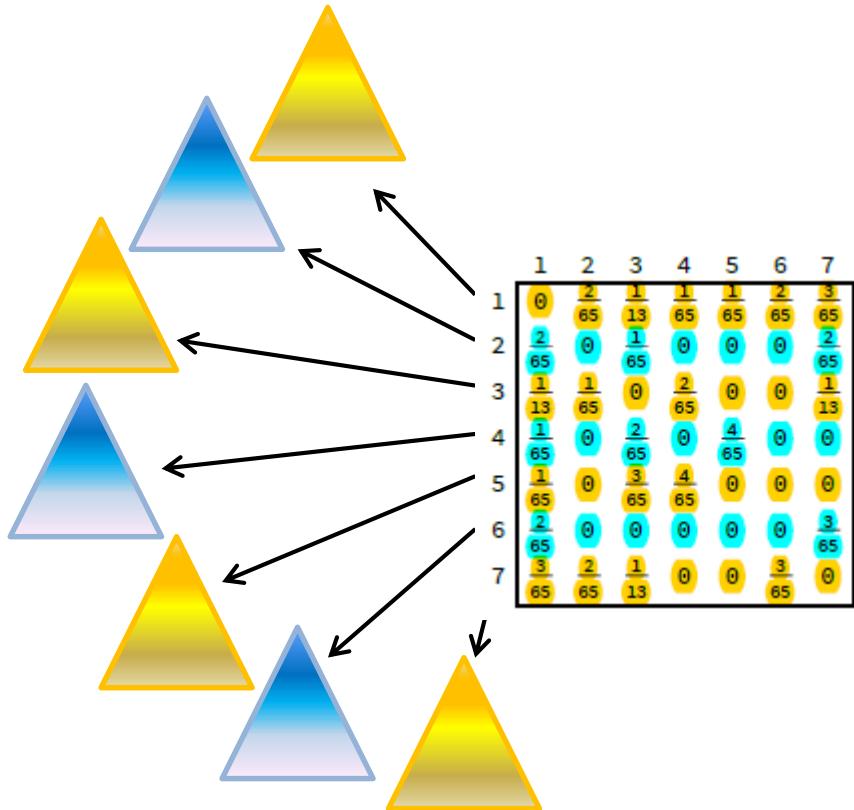
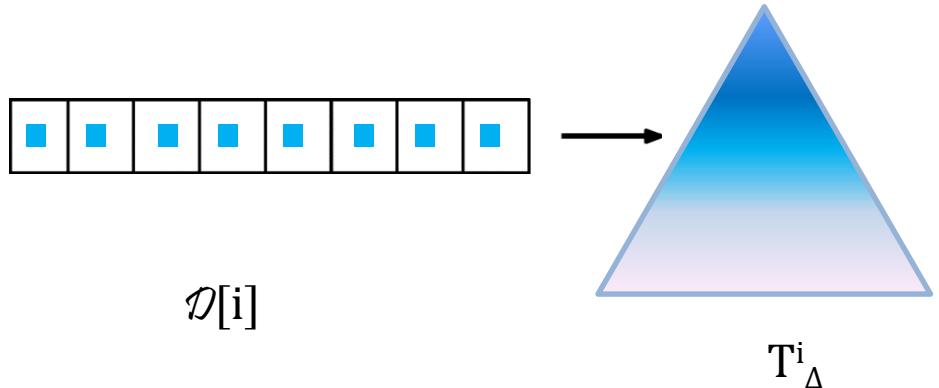
Ego-Trees Revisited

- Recall: ego-tree
 - optimal tree for a row (= given source)



Ego-Trees Revisited

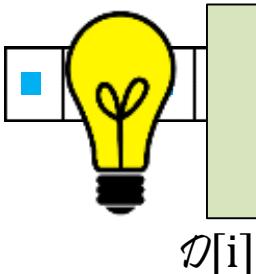
- Recall: ego-tree
 - optimal tree for a row (= given source)



Can we merge the trees **without distortion** and **keep degree low**?

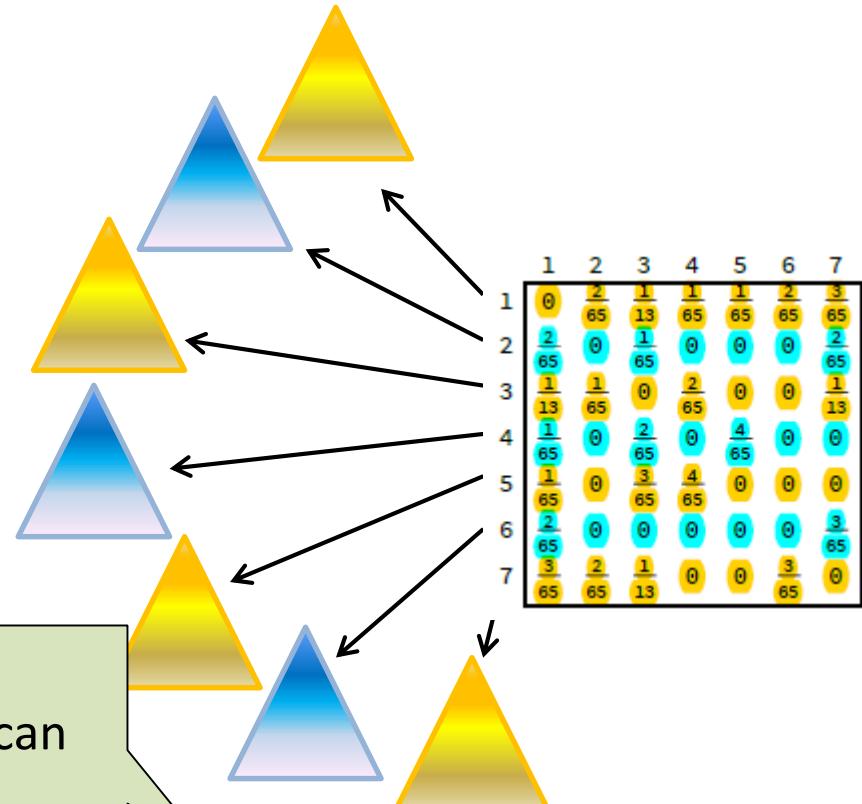
Ego-Trees Revisited

- Recall: ego-tree
 - optimal tree for a row (= given source)



For **sparse demands** yes:
enough *low-degree nodes* which can
serve as “*helper nodes*”!

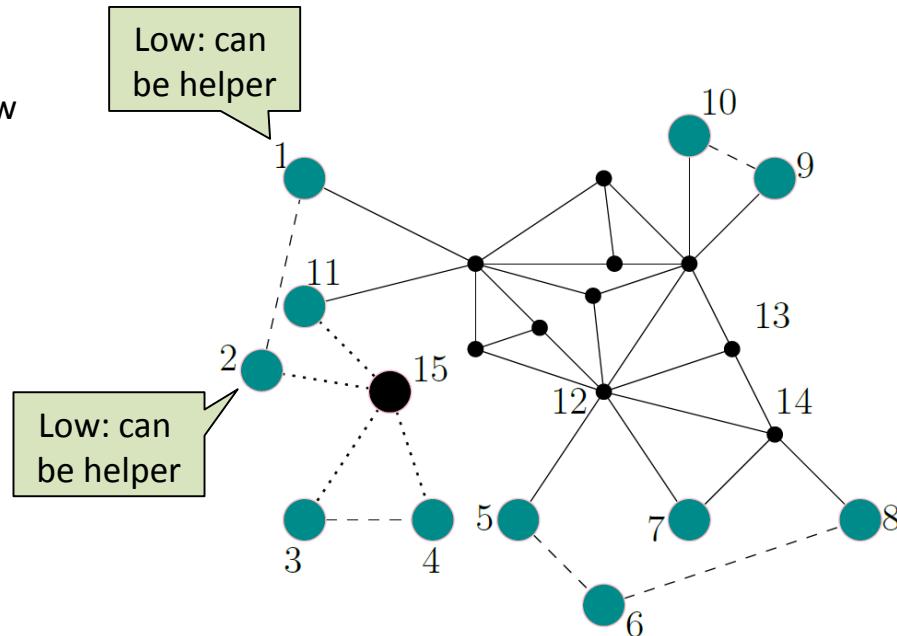
T_{Δ}^i



Can we merge the trees **without distortion** and **keep degree low**?

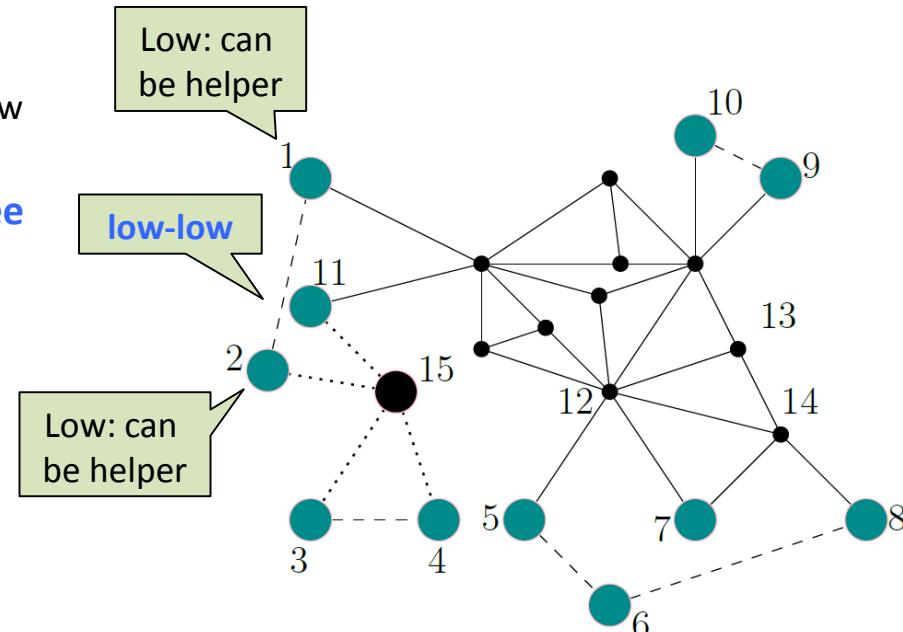
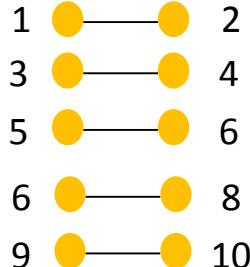
DAN for Sparse Demand

- Find **low** degree nodes
 - Half of the nodes of lowest degree: “below twice average degree”



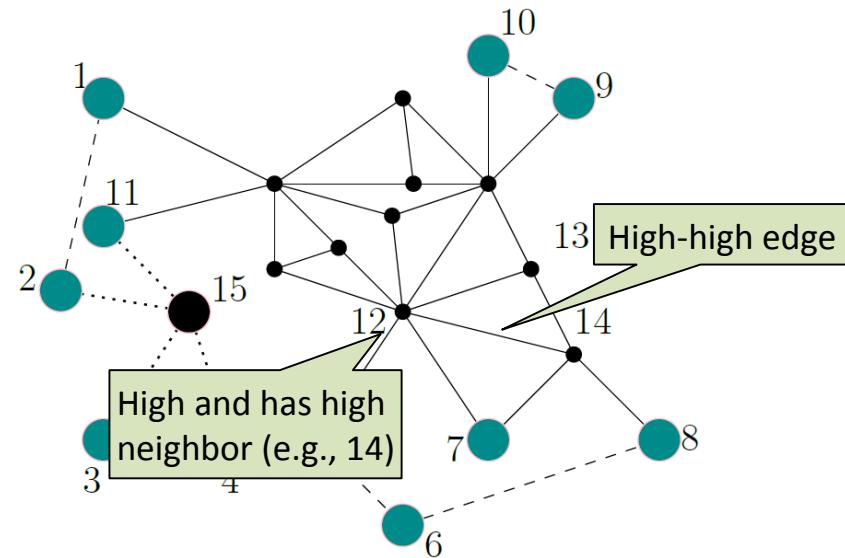
DAN for Sparse Demand

- Find **low** degree nodes
 - Half of the nodes of lowest degree: “below twice average degree”
- **Put** the **low-low** edges and the binary **tree** into DAN and remove from demand



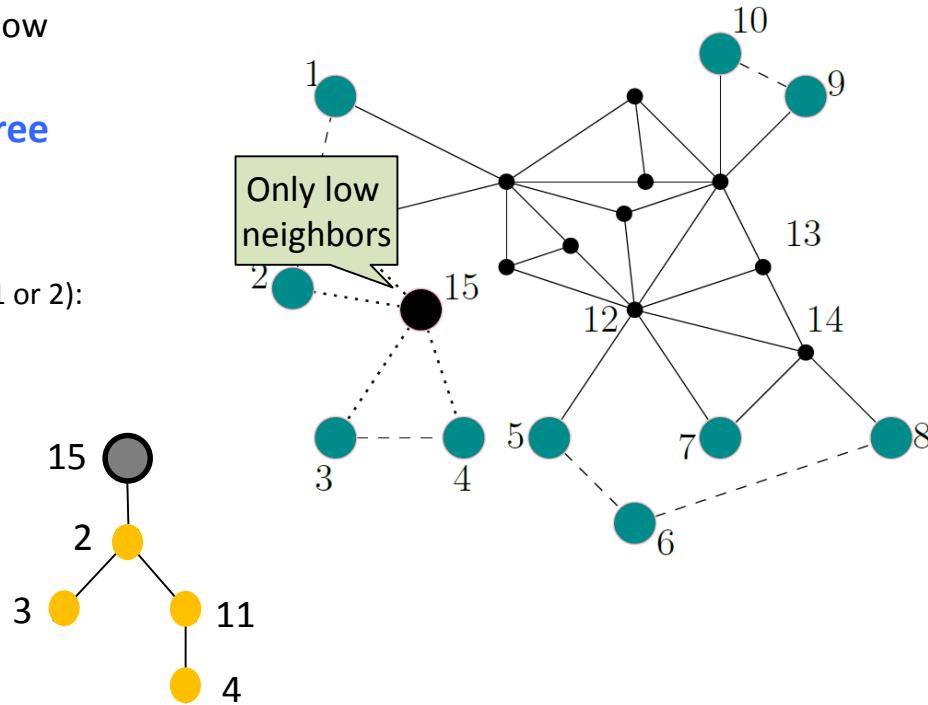
DAN for Sparse Demand

- Find **low** degree nodes
 - Half of the nodes of lowest degree: “below twice average degree”
 - Put the **low-low** edges and the binary **tree** into DAN and remove from demand
 - Mark **high-high** edges
 - Put (any) **low degree** nodes in between (e.g., 1 or 2): one is enough so distance increased by **+1**



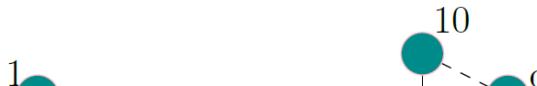
DAN for Sparse Demand

- Find **low** degree nodes
 - Half of the nodes of lowest degree: “below twice average degree”
- **Put** the **low-low** edges and the **binary tree** into DAN and remove from demand
- Mark **high-high** edges
 - Put (any) **low degree** nodes in between (e.g., 1 or 2): one is enough so distance increased by **+1**
- Now high degree nodes have only low degree neighbors: make **tree**
 - Create optimal **binary tree** with low degree neighbors



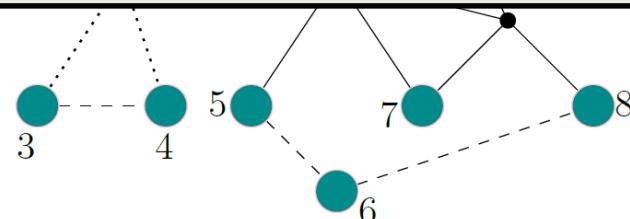
DAN for Sparse Demand

- Find **low** degree nodes
 - Half of the nodes of lowest degree: “below twice average degree”



Theorem [Asymptotic Optimality]: Helper node does not participate in many trees, so ***constant degree***, and ***constant distortion***.

- Now high degree nodes have only low degree neighbors: make **tree**
 - Create optimal **binary tree** with low degree neighbors



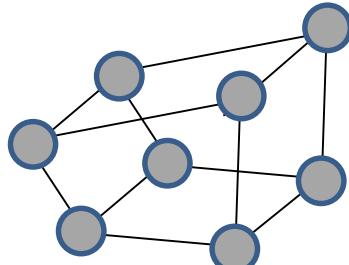
Remark: The Problem is Related To Spanners

- Sparse, distance-preserving (low-distortion) **spanners**
- But:
 - Spanners aim at low distortion among ***all pairs***; in our case, we are only interested in the ***local distortion***, 1-hop communication neighbors
 - We allow **auxiliary edges** (not a subgraph): similar to **geometric spanners**
 - We require ***constant degree***

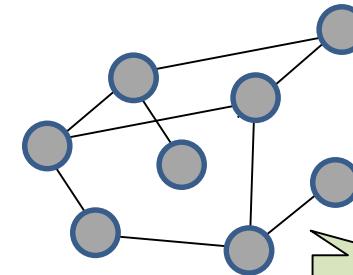
Yet: We can leverage the connection to spanners sometimes!

Theorem: If request distribution \mathcal{D} is **regular and uniform**, and if we can find a constant distortion, linear sized (i.e., **constant, sparse**) spanner for this request graph: then we can design a constant degree DAN providing an **optimal EPL** (i.e., $O(H(X|Y) + H(Y|X))$).

r-regular and uniform
demand:

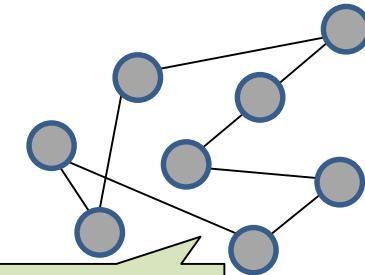


*Sparse, irregular
(constant) spanner:*



subgraph!

Constant degree optimal
DAN (EPL at most **log r**):



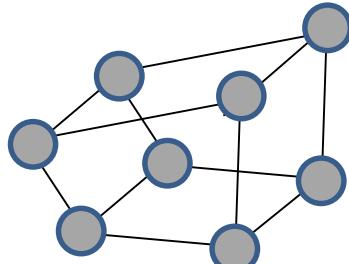
auxiliary edges

Yet: We can leverage the connection to spanners sometimes!

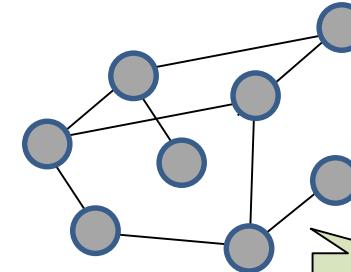
Theorem: If request distribution \mathcal{D} is **regular and uniform**, and if we can find a constant distortion, linear sized (i.e., **constant, sparse**) spanner for this request graph: then we can design a constant degree DAN providing

Optimal: in r -regular graphs, $\text{EPL}(\mathcal{Y}|\mathcal{X}) = \log r$.

r-regular and uniform
demand:

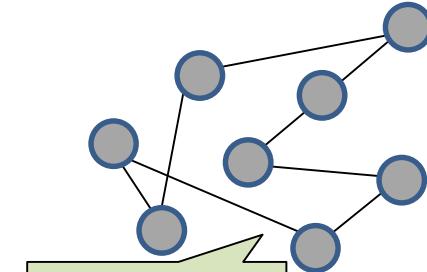


***Sparse, irregular
(constant) spanner:***



subgraph!

***Constant degree optimal
DAN (EPL at most $\log r$):***

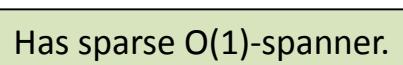


auxiliary edges

Proof Idea

- **Degree reduction** again, this time *from sparse spanner* (before: from sparse demand graph)

Corollaries

- Optimal DAN designs for
 - **Hypercubes** (with $n \log n$ edges)
 - **Chordal graphs**  Has sparse $O(1)$ -spanner.
 - Trivial: graphs with polynomial degree (**dense graphs**)
 - Graphs of **locally bounded doubling dimension**

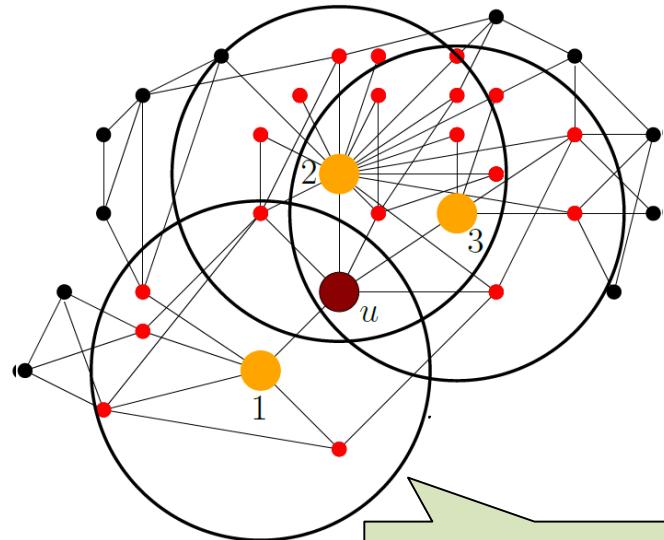
We also know
some more algos,
e.g., for BSTs.

Another Example: Demands of Locally-Bounded Doubling Dimension

- LDD: G_{\emptyset} has a **Locally-bounded Doubling Dimension** (LDD) iff all 2-hop neighbors are covered by 1-hop neighbors of just λ nodes
 - Note: care only about **2-neighborhood**

We only consider 2 hops!

- Formally, $B(u, 2) \subseteq \bigcup_{i=1}^{\lambda} B(v_i, 1)$
- Challenge: can be of **high degree**!



Nodes 1,2,3 cover 2-hop neighborhood of u .

DAN for Locally-Bounded Doubling Dimension

Lemma: There exists a sparse 9-(subgraph)spanner for LDD.

This *implies optimal DAN*: still focus on regular and uniform!

Def. (ε -net): A subset V' of V is a **ε -net** for a graph $G = (V, E)$ if

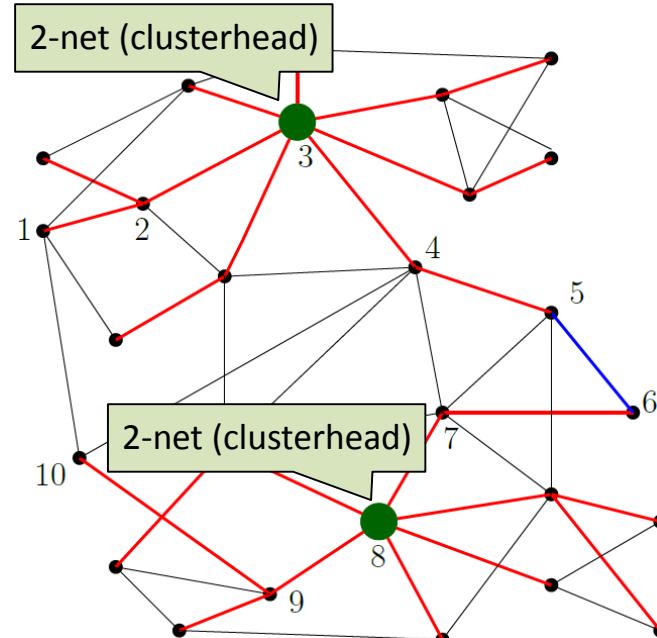
- V' sufficiently “**independent**”: for every $u, v \in V'$, $d_G(u, v) > \varepsilon$
- “**dominating**” V : for each $w \in V$, \exists at least one $u \in V'$ such that, $d_G(u, w) \leq \varepsilon$

9-Spanner for LDD (= optimal DAN)

Simple algorithm:

1. Find a 2-net

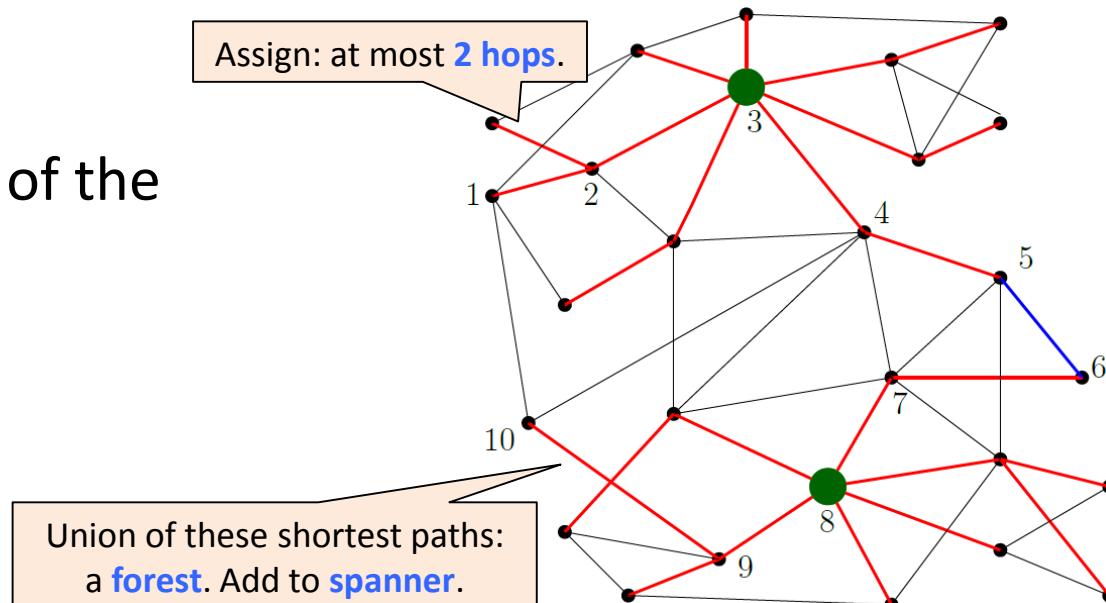
Easy: Select nodes into 2-net **one-by-one** in decreasing (remaining) degrees, **remove 2-neighborhood**. Iterate.



9-Spanner for LDD (= optimal DAN)

Simple algorithm:

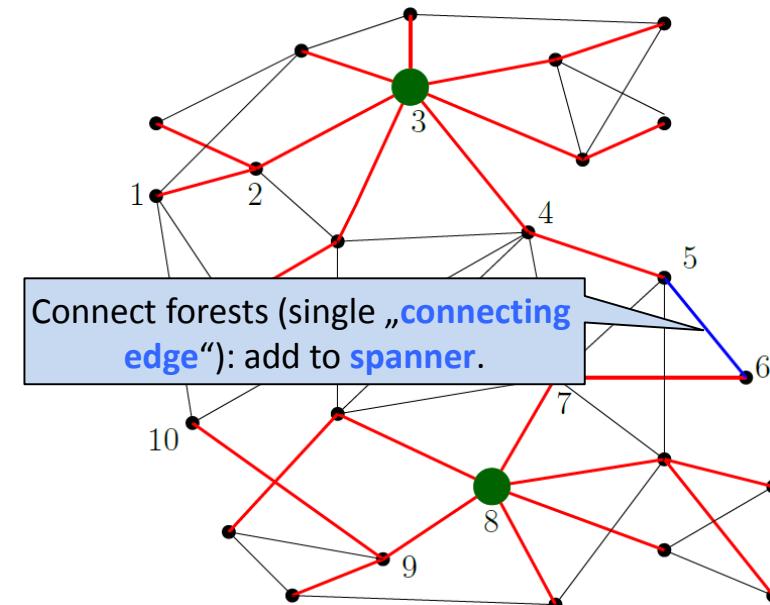
1. Find a 2-net
2. Add nodes to one of the closest 2-net nodes



9-Spanner for LDD (= optimal DAN)

Simple algorithm:

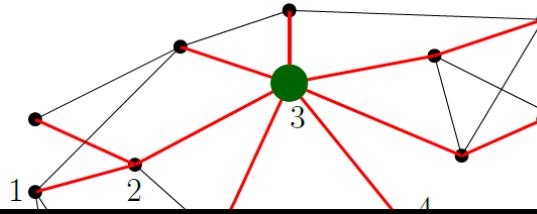
1. Find a 2-net
2. Add nodes to one of the closest 2-net nodes
3. Join two clusters if there are edges in between



9-Spanner for LDD (= optimal DAN)



Distortion 9: *Short detour* via
clusterheads: $u, ch(u), x, y, ch(v), v$



2. Add nodes to one of the
closest 2-net nodes

3. Join two clusters
edges in between



Sparse: Spanner only includes *forest* (sparse) plus
“connecting edges”: but since in *a locally doubling dimension graph* the number of cluster heads at
distance 5 is bounded, only a small number of
neighboring clusters will communicate.



Further Reading

Demand-Aware Network Designs of Bounded Degree

Chen Avin, Kaushik Mondal, and Stefan Schmid.

31st International Symposium on Distributed Computing (**DISC**),

Vienna, Austria, October 2017.

Roadmap

- Motivation: Demand-Aware Networks
- Principles of Static Demand-Aware Network Designs
- **Principles of Dynamic Demand-Aware Network Designs**
- Principles of Decentralized Approaches



Objectives and Metrics for Dynamic DANs, i.e. SANs?

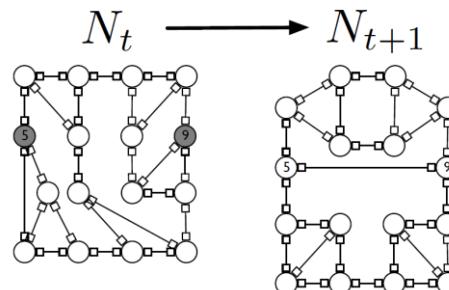
Input for Dynamic DANs

A **sequence** $\sigma = (u_1, v_1), (u_2, v_2), (u_3, v_3) \dots$

chosen arbitrarily

Chosen i.i.d. from initially
unknown fixed distribution

A Cost-Benefit Tradeoff



Basic question:

How often to reconfigure?

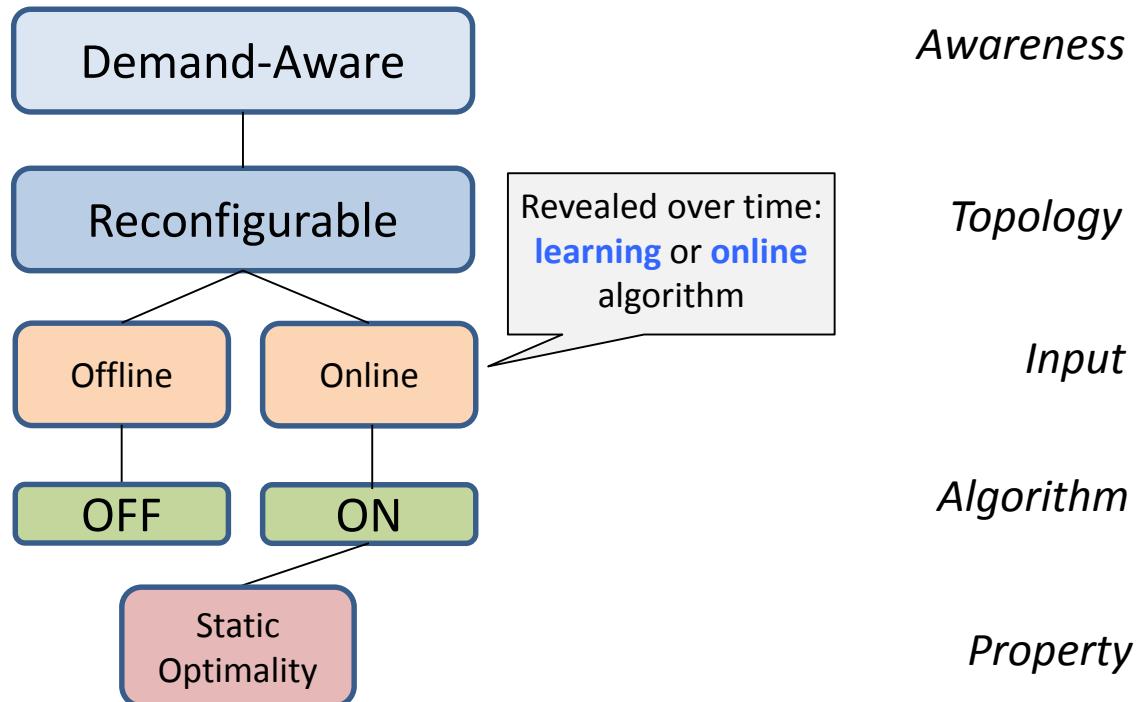
Short routes
High reconfiguration cost



Low reconfiguration cost
Long routes

A Taxonomy: Reconfigurable Networks

Static Optimality:
“Not worse than static
which knows demand
ahead of time!”
 $\rho = \text{Cost(ON)}/\text{Cost(STAT*)}$
is constant.



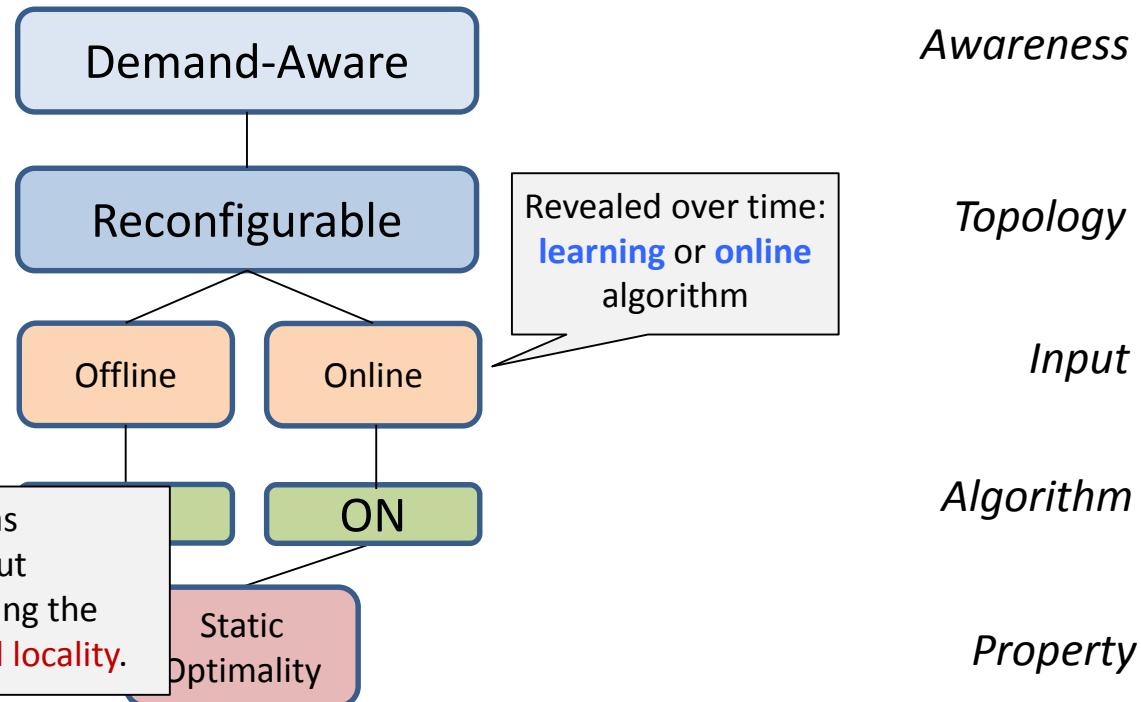
A Taxonomy: Reconfigurable Networks

Static Optimality:

"Not worse than static
which knows demand
ahead of time!"

$\rho = \text{Cost(ON)}/\text{Cost(STAT*)}$
is constant.

Note: may be $<<1$. ON has
advantage of adjusting, but
the **disadvantage** of not knowing the
workload. E.g. if much **temporal locality**.



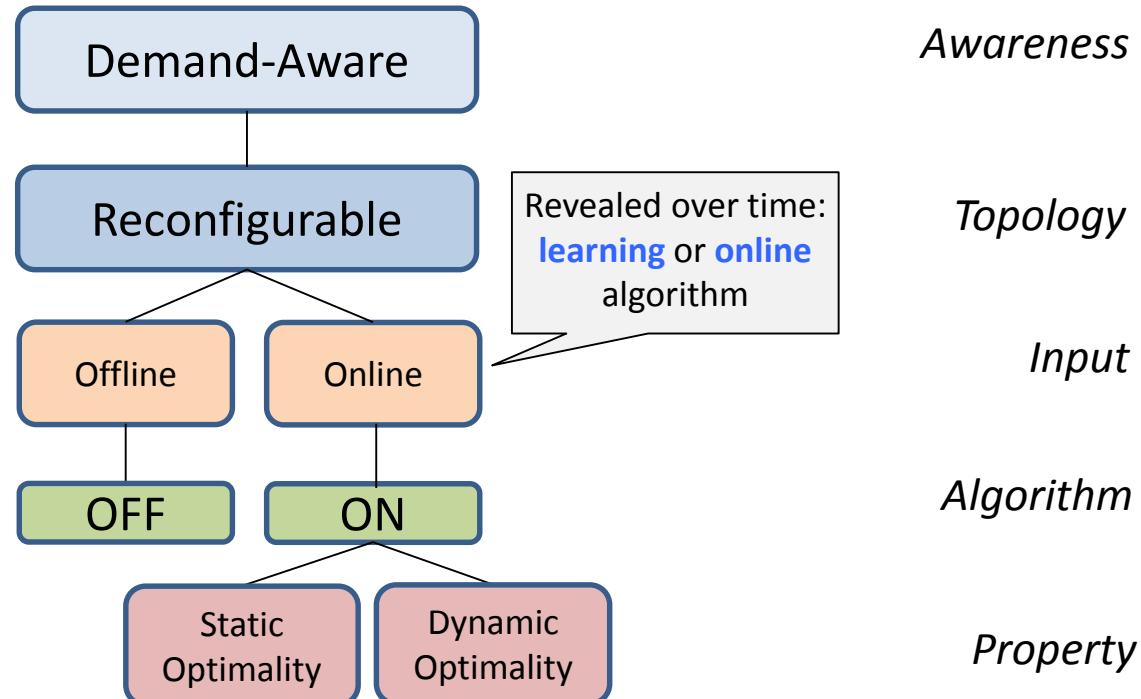
A Taxonomy: Reconfigurable Networks

Dynamic Optimality:

No worse than an offline algorithm which ***knows the sequence!***"

$\rho = \text{Cost(ON)}/\text{Cost(OFF*)}$
is constant.

Always ≥ 1 .



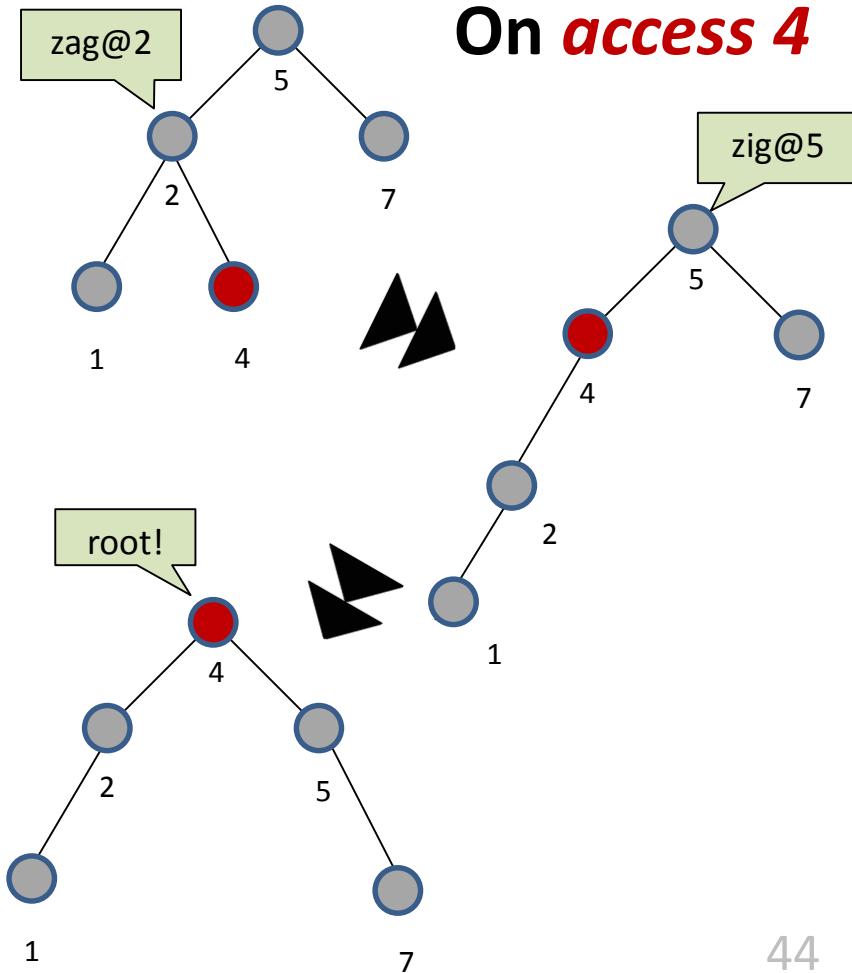
How to Design SANs?



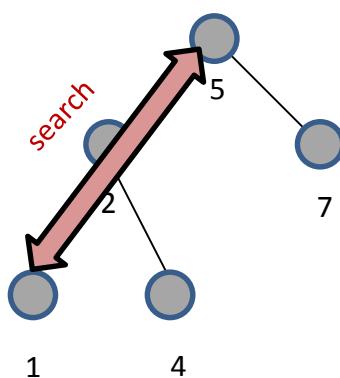
Inspiration from **self-adjusting
datastructures** again!

Recall: Splay Tree

- A Binary Search Tree (**BST**)
- Inspired by “**move-to-front**”: move **to root!**
- Self-adjustment: **zig**, **zigzag**, **zigzag**
 - Maintains **search property**
- Many nice properties
 - **Static optimality**, **working set**, (static,dynamic) **fingers**, ...

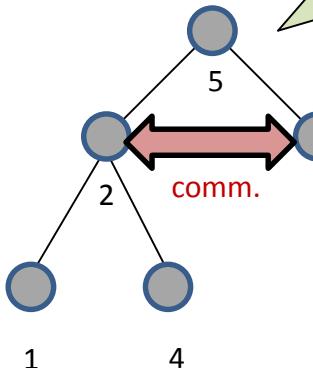


A Simple Idea: Generalize Splay Tree To *SplayNet*



Splay Tree

vs

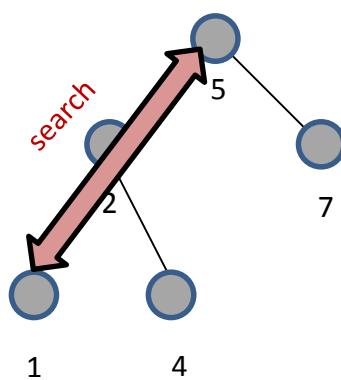


SplayNet

BST is nice for networks:
local (greedy) search!

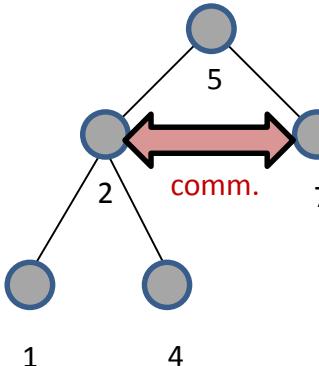
A Simple Idea: Generalize Splay Tree To *SplayNet*

But how?



Splay Tree

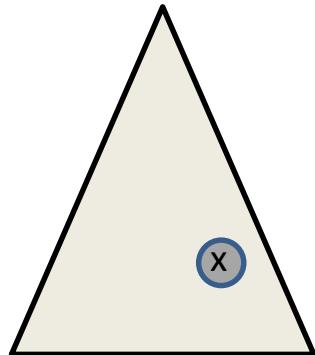
vs



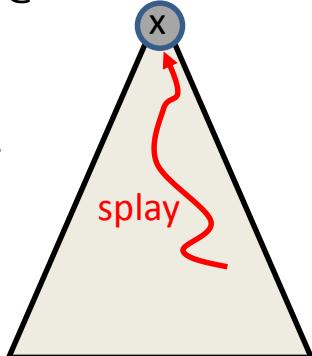
SplayNet

SplayNet: A Simple Idea

@t: access x

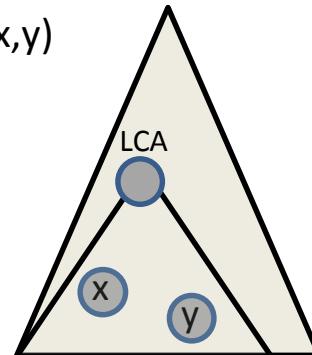


@t+1

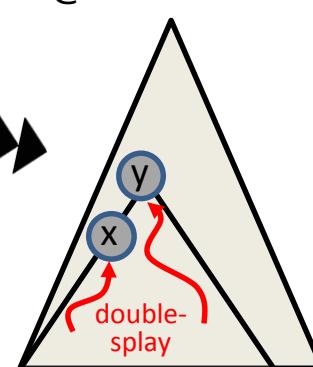


Splay Tree

@t: comm
(x,y)

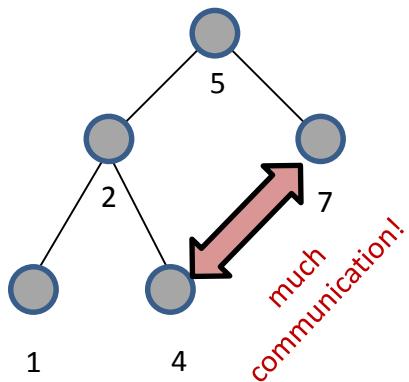


@t+1



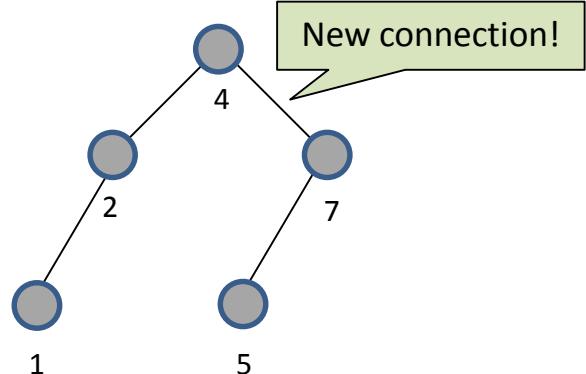
SplayNet

Example



$t=1$

adjust



$t=2$

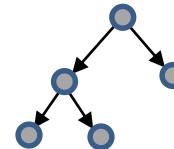
Challenges: How to minimize reconfigurations?

How to keep network locally routable?

Properties of SplayNets

- **Statically optimal** if demand comes from a *product distribution*
 - Product distribution: entropy equals conditional entropy, i.e., $H(X)+H(Y)=H(X|Y)+H(Y|X)$
- Converges to optimal static topology in
 - **Multicast scenario**: requests come from a BST as well
 - **Cluster scenario**: communication only *within* interval
 - **Laminated scenario** : communication is „non-crossing matching“

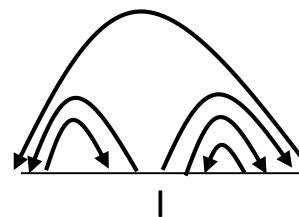
Multicast Scenario



Cluster Scenario



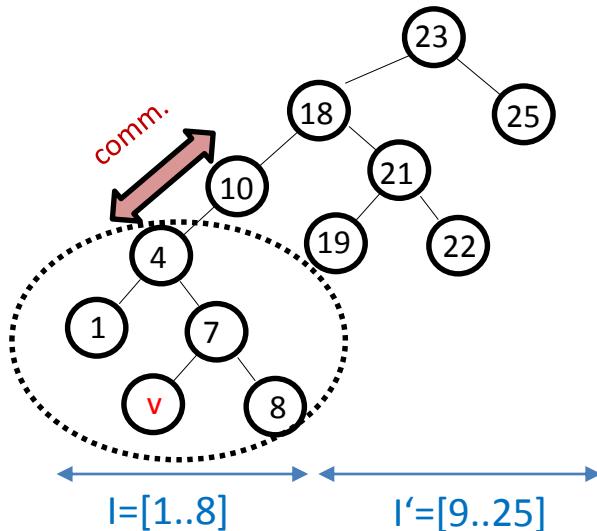
Laminated Scenario



Remark: Static SplayNet

Theorem: Optimal static SplayNet can be computed in polynomial-time (**dynamic programming**)

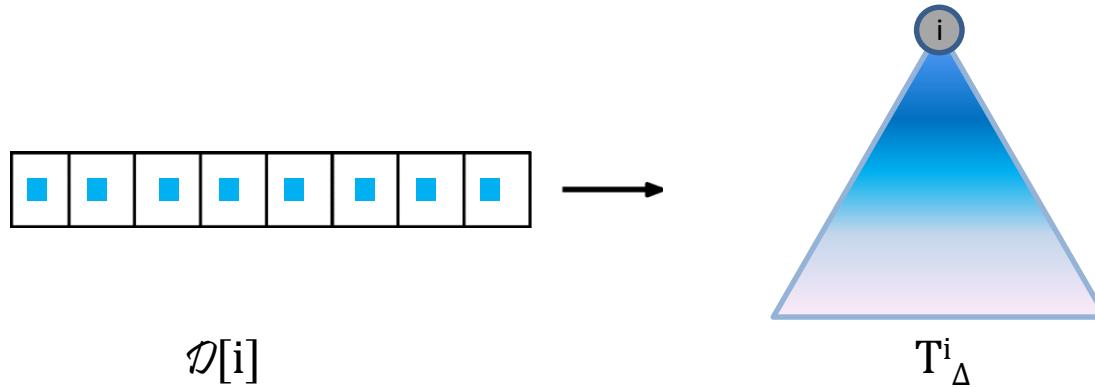
- Unlike unordered tree?



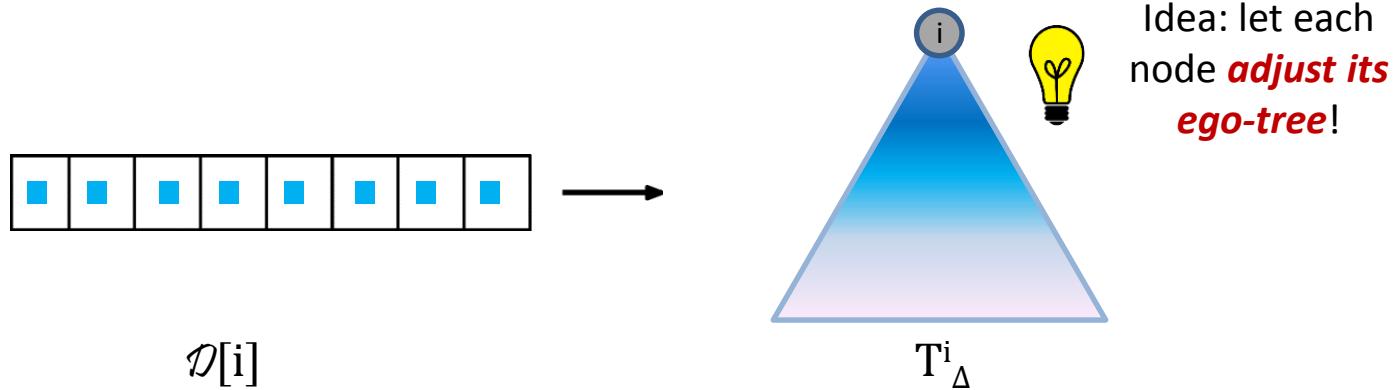
Further Reading

SplayNet: Towards Locally Self-Adjusting Networks
Stefan Schmid, Chen Avin, Christian Scheideler, Michael Borokhovich, Bernhard Haeupler, and Zvi Lotker.
IEEE/ACM Transactions on Networking (**TON**), Volume 24, Issue 3, 2016.

Better Idea: Back to Ego-Trees!

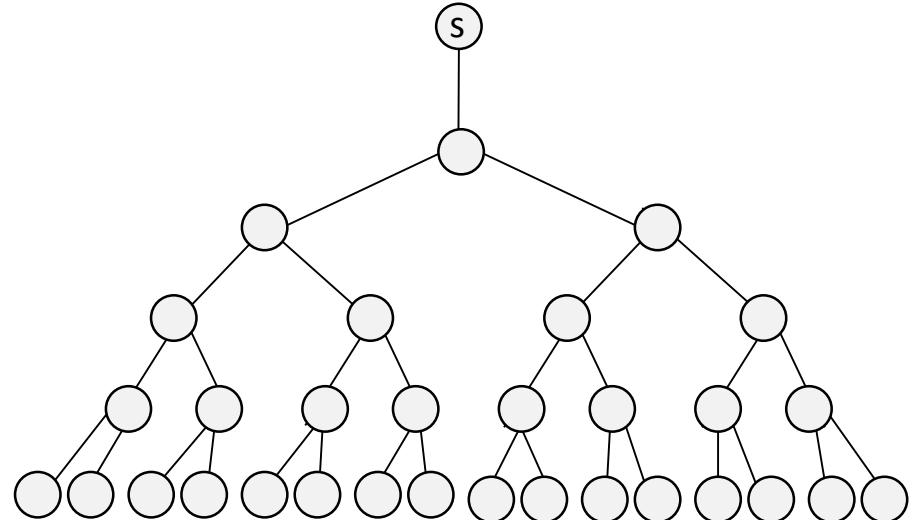


Better Idea: Back to Ego-Trees!



A Balanced Self-Adjusting Tree: Push-Down Tree

- **Push-down tree:** a self-adjusting complete tree
- ***Dynamically optimal***
- Not ordered: requires **a map**

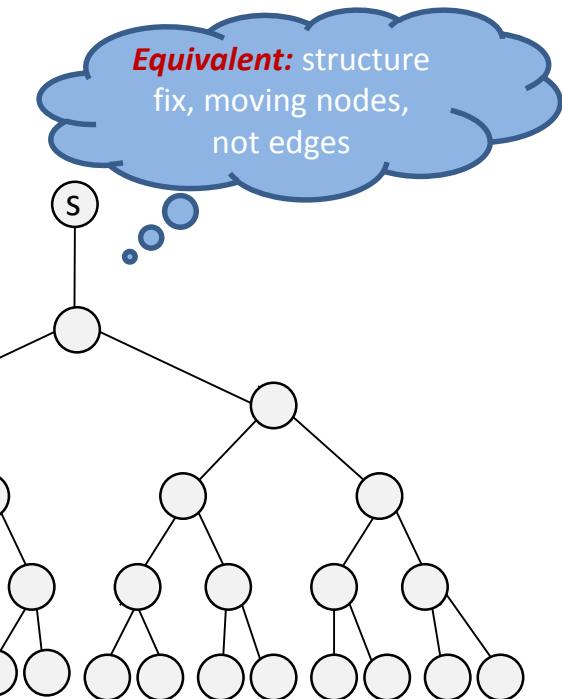


A Balanced Self-Adjusting Tree:

Push-Down Tree



- **Push-down tree:** a self-adjusting complete tree
- ***Dynamically optimal***
- Not ordered: requires **a map**



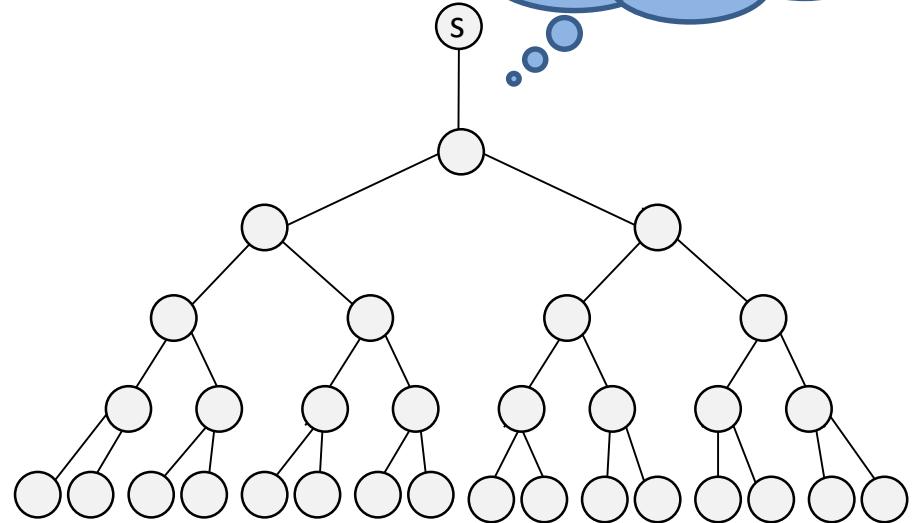
A Balanced Self-Adjusting Tree:

Push-Down Tree



- **Push-down tree:** a self-adjusting complete tree
- **Dynamically optimal**
- Not ordered: requires **a map**

A blue cloud-shaped graphic containing the text "Equivalent: structure fix, moving nodes, not edges" in red and black font. Below the cloud are three small blue dots.

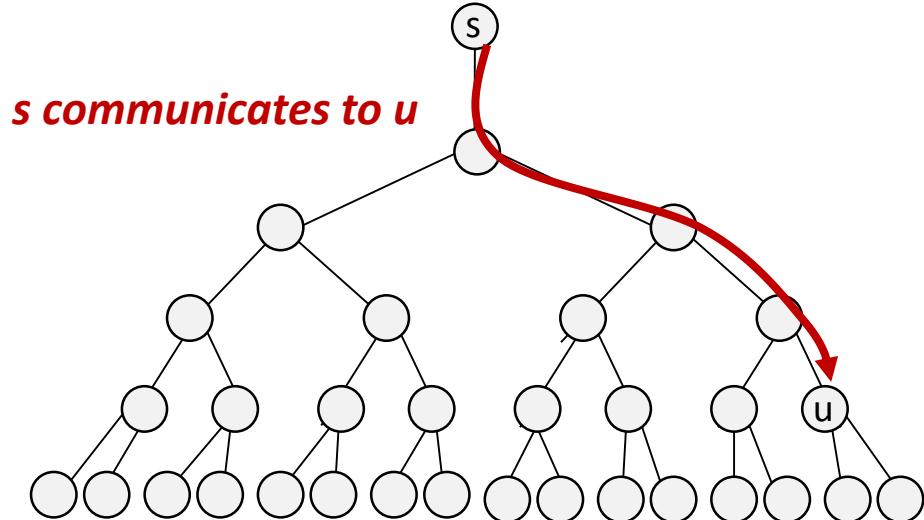


A useful dynamic property: **Most-Recently Used (MRU)**!

Similar to **Working Set Property**: more recent communication Partners closer to source.

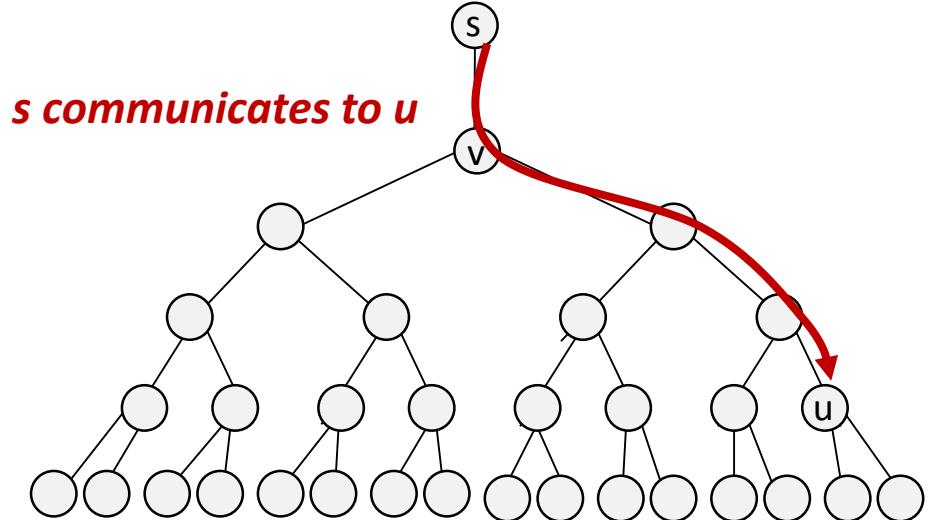
A Balanced Self-Adjusting Tree: Push-Down Tree

- **Push-down tree:** a self-adjusting complete tree
- *Dynamically optimal*
- Not ordered: requires **a map**



A Balanced Self-Adjusting Tree: Push-Down Tree

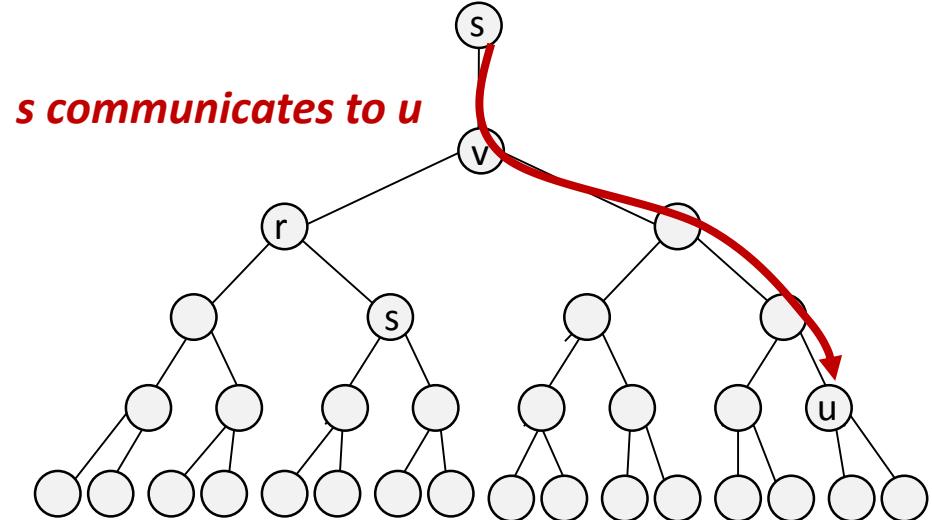
- **Push-down tree:** a self-adjusting complete tree
- **Dynamically optimal**
- Not ordered: requires **a map**



Strict MRU requires: move u to root! But how? Cannot swap with v: v no longer MRU!

A Balanced Self-Adjusting Tree: Push-Down Tree

- **Push-down tree:** a self-adjusting complete tree
- **Dynamically optimal**
- Not ordered: requires **a map**



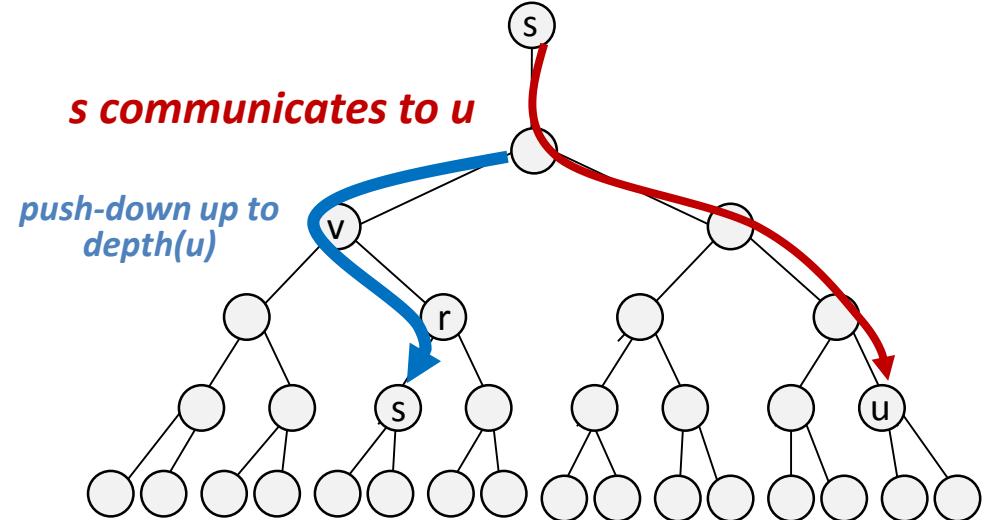
Strict MRU requires: move u to root! But how? Cannot swap with v: v no longer MRU!



Idea: Push v down, in a balanced manner, up to depth(u): left-right-left-right („rotate-push“)

A Balanced Self-Adjusting Tree: Push-Down Tree

- **Push-down tree:** a self-adjusting complete tree
- **Dynamically optimal**
- Not ordered: requires **a map**



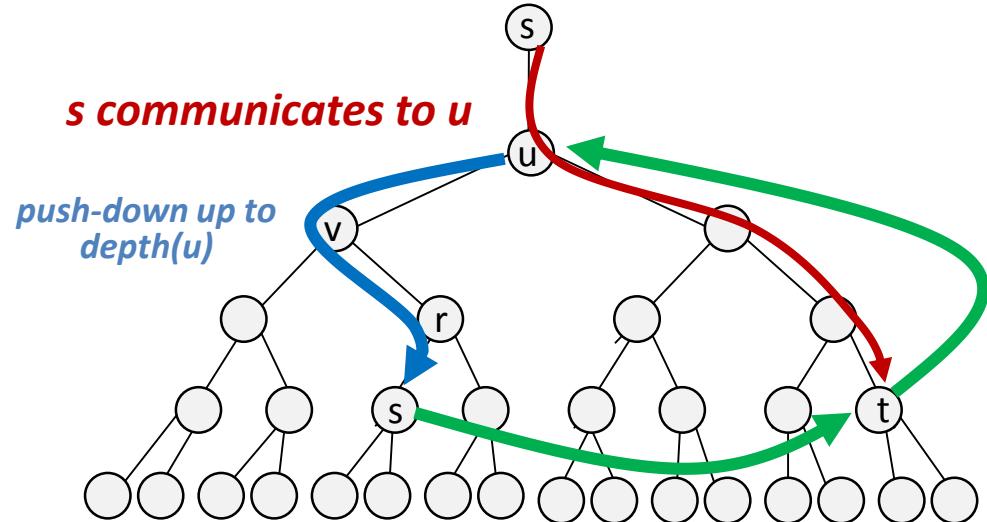
Strict MRU requires: move u to root! But how? Cannot swap with v: v no longer MRU!



Idea: Push v down, in a balanced manner, up to depth(u): left-right-left-right („rotate-push“)

A Balanced Self-Adjusting Tree: Push-Down Tree

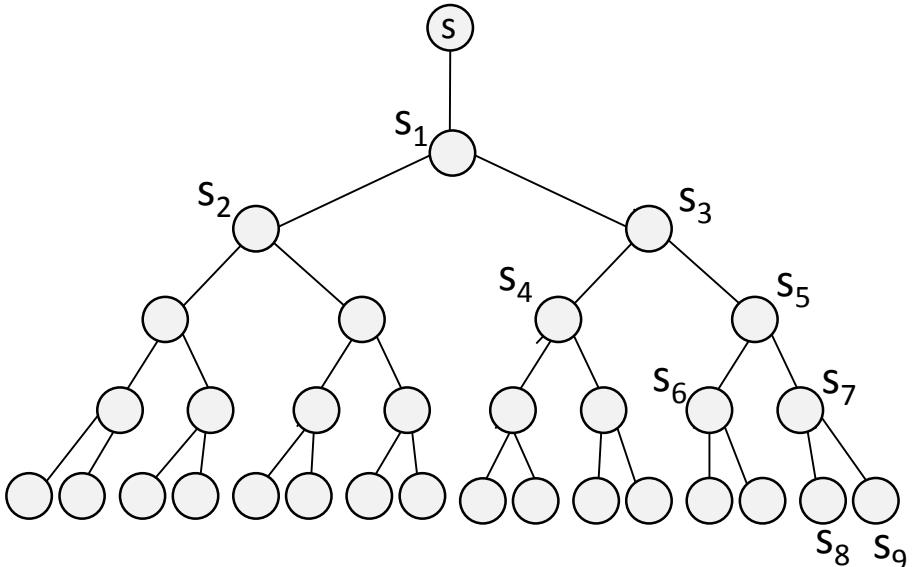
- **Push-down tree:** a self-adjusting complete tree
- **Dynamically optimal**
- Not ordered: requires **a map**



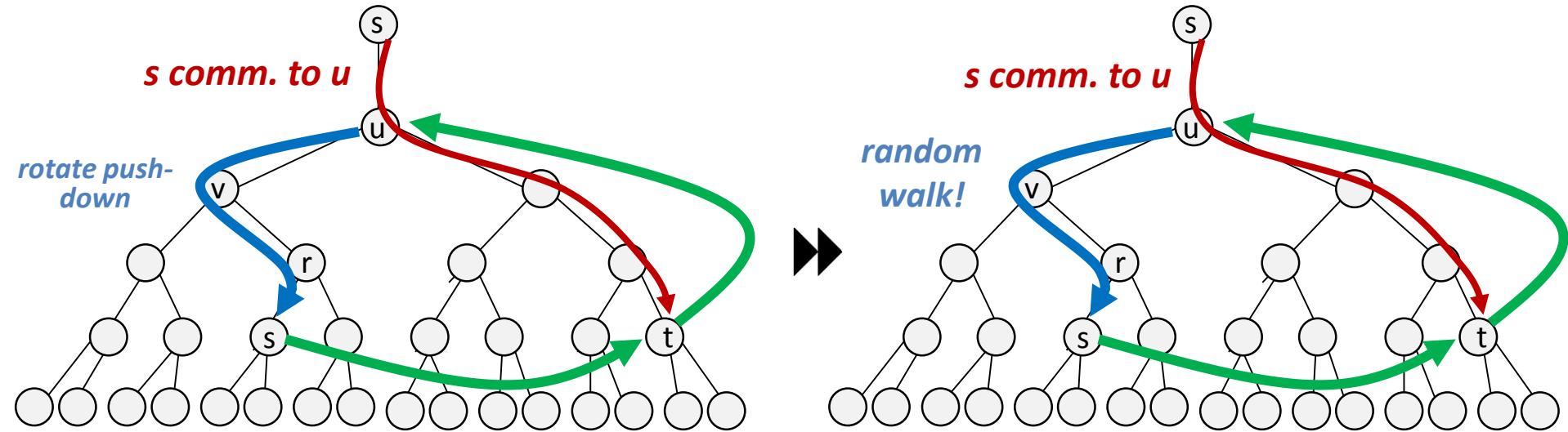
*Then: promote u to available root, and
t to u: at original depth!*

Remarks

- Unfortunately, alternating push-down does ***not maintain MRU*** (working set) property
- Tree can ***degrade***, e.g.: sequence of requests from level 4,1,2,1,3,1,4,1



Solution: Random Walk



*At least maintains approximate
working set / MRU!*

Further Reading

Push-Down Trees: Optimal Self-Adjusting Complete Trees

Chen Avin, Kaushik Mondal, and Stefan Schmid.

ArXiv Technical Report, July 2018.

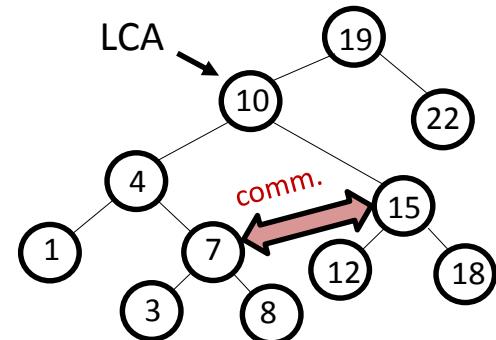
Roadmap

- Motivation: Demand-Aware Networks
- Principles of Static Demand-Aware Network Designs
- Principles of Dynamic Demand-Aware Network Designs
- **Principles of Decentralized Approaches**



A “Simple” Decentralized Solution: Distributed SplayNet (*DiSplayNet*)

- SplayNet attractive: ordered BST supports **local routing**
 - Nodes **maintain three ranges**: interval of left subtree, right subtree, upward
- If communicate (frequently): **double-splay** toward LCA
- Challenge: **concurrency**!
 - Access Lemma of splay trees no longer works: **potential function** does not „**telescope**“ anymore: a concurrently rising node may push down another rising node again

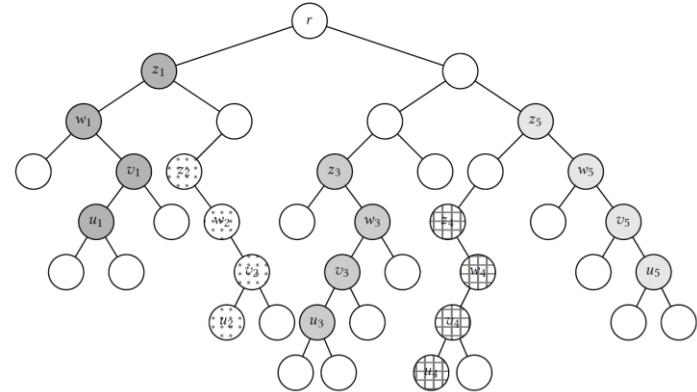


SplayNet

DiSplayNet: Challenges

- DiSplayNet: Rotations (zig,zigzag,zigzag) are **concurrent**
- To avoid conflict: distributed computation of **independent clusters**
- Still challenging:

	1	2	3	4	5	6	7	8	...	$i - 6$	$i - 5$	$i - 4$	$i - 3$	$i - 2$	$i - 1$	i
σ_1	✓	✓	✓	✓	-	-	-	-	...	-	-	-	-	-	-	-
σ_2	-	X	X	X	✓	✓	✓	-	...	-	-	-	-	-	-	-
...
σ_{m-1}	-	-	-	-	-	-	-	-	...	✓	✓	-	-	-	-	-
σ_m	-	-	-	-	-	-	-	-	...	X	X	✓	✓	✓	✓	-



	1	2	3	...	i	$i + 1$	$i + 2$	$i + 3$	$i + 4$	$i + 5$	$i + 6$...	j	...	k
s_1	✓	✓	✓	...	✓	✓	✓	✓	✓	✓	...	-	✓	...	-
d_1	✓	✓	✓	...	✓	✓	✓	✓	✓	✓	...	-	✓	...	-
s_2	-	✓	✓	...	✓	✓	✓	✓	-	-	...	-	-	...	-
d_2	-	✓	✓	...	✓	✓	X	✓	-	-	...	-	-	...	-
s_3	-	-	✓	...	X	X	X	✓	X	X	...	✓	...	-	-
d_3	-	-	✓	...	X	X	X	X	X	X	...	✓	...	-	-

Sequential SplayNet: requests **one after another**

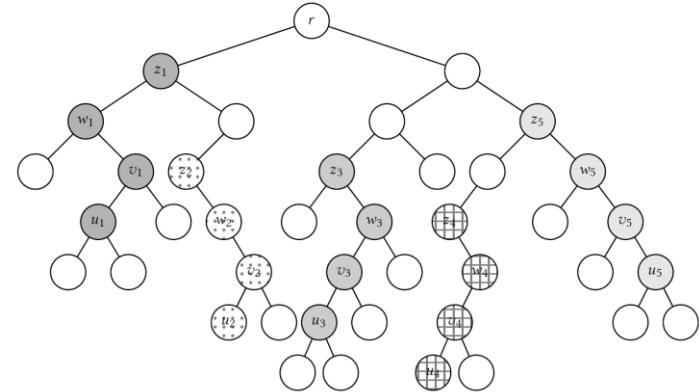
DiSplayNet: Analysis more challenging: potential function sum no longer **telescopic**. One request can “push-down” another.

DiSplayNet: Challenges

- DiSplayNet: Rotations (zig,zigzag,zigzag) are **concurrent**
- To avoid conflict: distributed computation of **independent clusters**
- Still challenging: **Telescopic: max potential drop**

	1	2	3	4	5	6	7	8	...	$i - 6$	$i - 5$	$i - 4$	$i - 3$	$i - 2$	$i - 1$	i
σ_1	✓	✓	✓	✓				-	...	-	-	-	-	-	-	-
σ_2	-	X	X	X	✓	✓	✓	-	...	-	-	-	-	-	-	-
...
σ_{m-1}	-	-	-	-	-	-	-	-	...	✓	✓				-	-
σ_m	-	-	-	-	-	-	-	-	...	X	X	✓	✓	✓	✓	-

Sequential SplayNet: requests **one after another**



	1	2	3	...	i	$i + 1$	$i + 2$	$i + 3$	$i + 4$	$i + 5$	$i + 6$...	j	...	k
s_1	✓	✓	✓	...	✓	✓	✓	✓	✓	✓	...	-	✓	...	-
d_1	✓	✓	✓	...	✓	✓	✓	✓	✓	✓	...	-	✓	...	-
s_2	-	✓	✓	...	✓	✓	✓	✓	-	-	...	-	-	...	-
d_2	-	✓	✓	...	✓	✓	X	✓	-	-	...	-	-	...	-
s_3	-	-	✓	...	X	X	X	✓	X	X	...	✓	...	-	-
d_3	-	-	✓	...	X	X	X	X	X	X	...	✓	...	-	-

DiSplayNet: Analysis more challenging: potential function sum no longer **telescopic**. One request can “push-down” another.

Further Reading

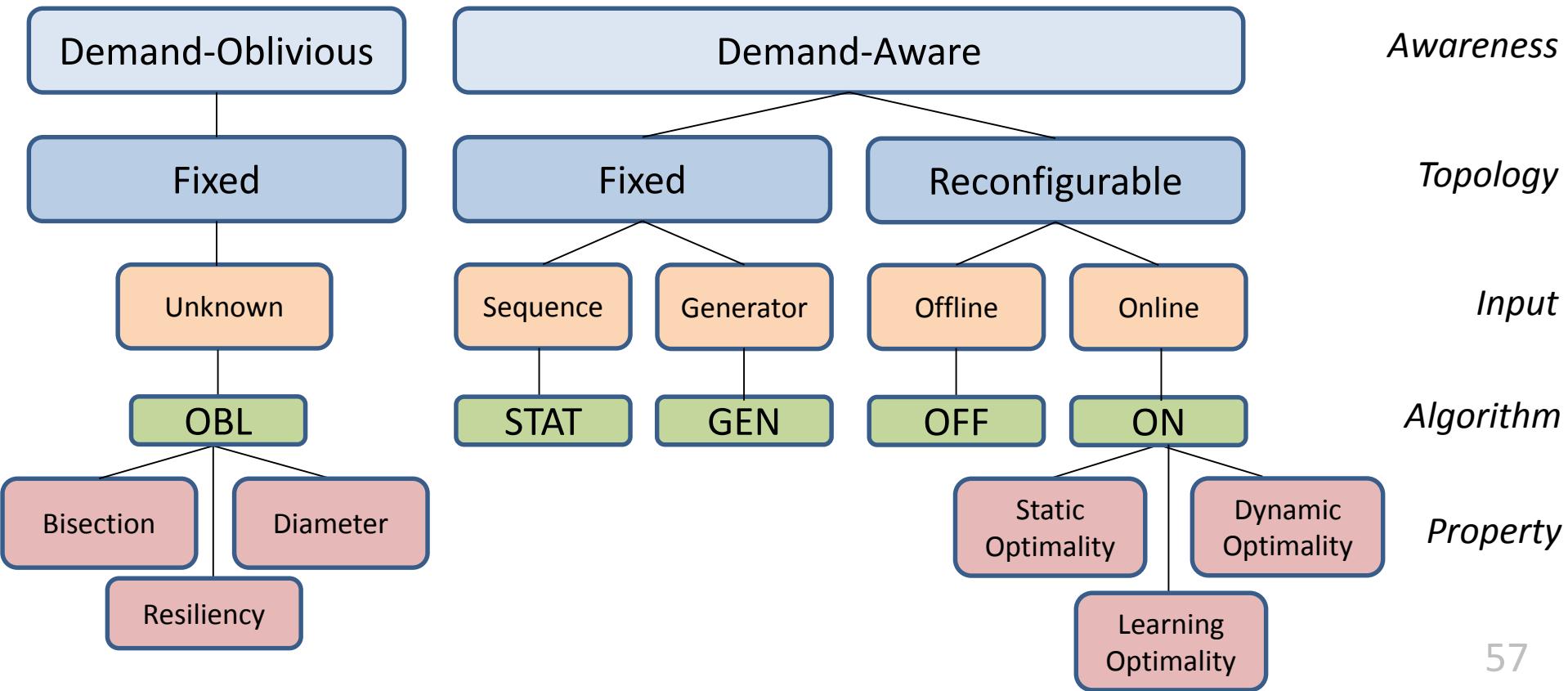
Brief Announcement: Distributed SplayNets

Bruna Peres, Olga Goussevskaia, Stefan Schmid, and Chen Avin.

31st International Symposium on Distributed Computing (**DISC**), Vienna,
Austria, October 2017.

Uncharted Landscape!

Toward Demand-Aware Networking: A Theory for
Self-Adjusting Networks. SIGCOMM CCR, 2018.



Conclusion

- Reconfigurable switches: *Yoga for Networks?*
- New metrics needed: e.g., entropy?
- New algorithms needed: static, offline and online!
- *Let's chat!*

*Thank you!
Question?*

Further Reading

[Toward Demand-Aware Networking: A Theory for Self-Adjusting Networks](#)

Chen Avin and Stefan Schmid.

SIGCOMM CCR, October 2018.

[Demand-Aware Network Designs of Bounded Degree](#)

Chen Avin, Kaushik Mondal, and Stefan Schmid.

31st International Symposium on Distributed Computing (**DISC**), Vienna, Austria, October 2017.

[Push-Down Trees: Optimal Self-Adjusting Complete Trees](#)

Chen Avin, Kaushik Mondal, and Stefan Schmid.

ArXiv Technical Report, July 2018.

[Online Balanced Repartitioning](#)

Chen Avin, Andreas Loukas, Maciej Pacut, and Stefan Schmid.

30th International Symposium on Distributed Computing (**DISC**), Paris, France, September 2016.

[rDAN: Toward Robust Demand-Aware Network Designs](#)

Chen Avin, Alexandr Hercules, Andreas Loukas, and Stefan Schmid.

Information Processing Letters (**IPL**), Elsevier, 2018.

[SplayNet: Towards Locally Self-Adjusting Networks](#)

Stefan Schmid, Chen Avin, Christian Scheideler, Michael Borokhovich, Bernhard Haeupler, and Zvi Lotker.

IEEE/ACM Transactions on Networking (**TON**), Volume 24, Issue 3, 2016. Early version: IEEE **IPDPS** 2013.

[Characterizing the Algorithmic Complexity of Reconfigurable Data Center Architectures](#)

Klaus-Tycho Foerster, Monia Ghobadi, and Stefan Schmid.

ACM/IEEE Symposium on Architectures for Networking and Communications Systems (**ANCS**), Ithaca, New York, USA, July 2018.

[Charting the Complexity Landscape of Virtual Network Embeddings](#)

Matthias Rost and Stefan Schmid. **IFIP Networking**, Zurich, Switzerland, May 2018.